

# DATA MINING 2

## Ethics Principles: Explainability

---

Riccardo Guidotti

a.a. 2019/2020

# Definitions

---

**explanation** | ɛksplə'neɪʃ(ə)n |

noun

a statement or account that makes something clear: *the birth rate is central to any explanation of population trends.*

**interpret** | ɪn'təːprɪt |

verb (**interprets, interpreting, interpreted**) [*with object*]

1 explain the meaning of (information or actions): *the evidence is difficult to interpret.*

# What is “Explainable AI” ?

---

- **Explainable-AI** explores and investigates methods to produce or complement **AI models** to make **accessible and interpretable** the internal logic and the outcome of the algorithms, making such process **understandable by humans**.
- **Explicability**, understood as incorporating both **intelligibility** (*“how does it work?”*) for non-experts, e.g., patients or business customers, and for experts, e.g., product designers or engineers) and **accountability** (*“who is responsible for”*).
- 5 core principles for ethical AI:
  - beneficence, non-maleficence, autonomy, and justice
  - a new principle is needed in addition: **explicability**

# Motivating Examples

- Criminal Justice
  - People wrongly denied
  - Recidivism prediction
  - Unfair Police dispatch
- Finance:
  - Credit scoring, loan approval
  - Insurance quotes
- Healthcare
  - AI as 3<sup>rd</sup>-party actor in physician - patient relationship
  - Learning must be done with available data: cannot randomize cares given to patients!
  - Must validate models before use.

Opinion

The New York Times

OP-ED CONTRIBUTOR

## When a Computer Program Keeps You in Jail

The Big Read **Artificial intelligence**

+ Add to myFT

### Insurance: Robots learn the business of covering risk



**Stanford**  
MEDICINE

News Center

Email

Tweet

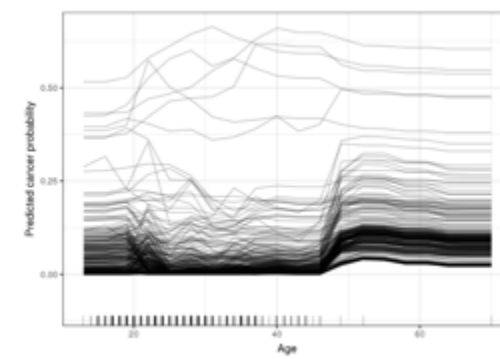
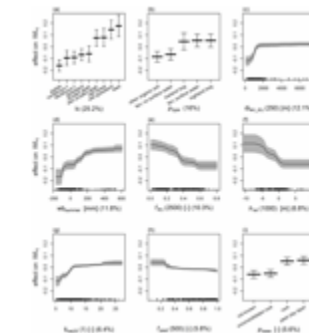
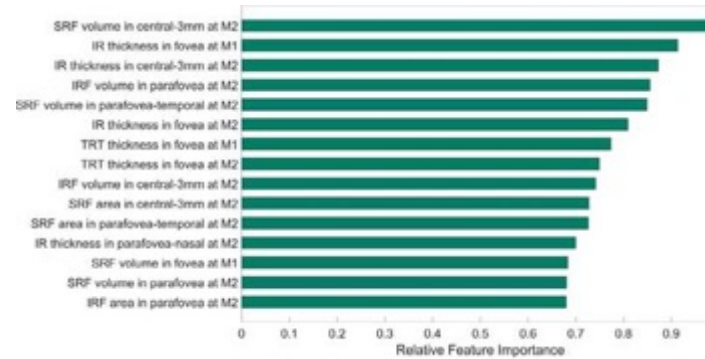
Researchers say use of artificial intelligence in medicine raises ethical questions

In a perspective piece, Stanford researchers discuss the ethical implications of using machine-learning tools in making health care decisions for patients.

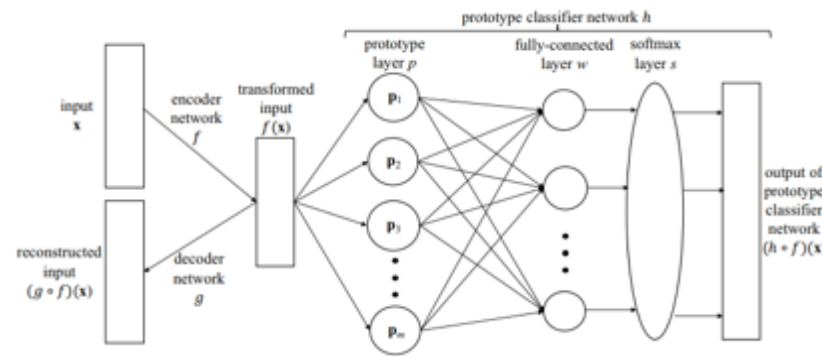


# Explanation in different AI fields

- Machine Learning

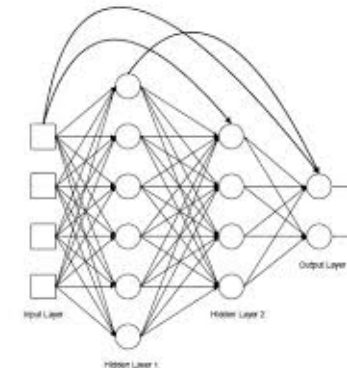


(a) Feature Importance, Partial Dependence Plot, Individual Conditional Expectation



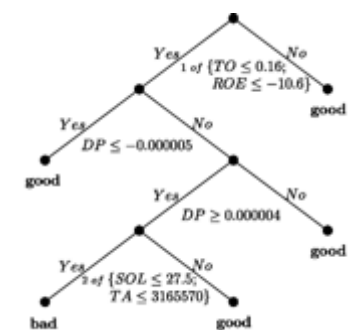
Auto-encoder

Oscar Li, Hao Liu, Chaofan Chen, Cynthia Rudin: Deep Learning for Case-Based Reasoning Through Prototypes: A Neural Network That Explains Its Predictions. AAAI 2018: 3530-3537



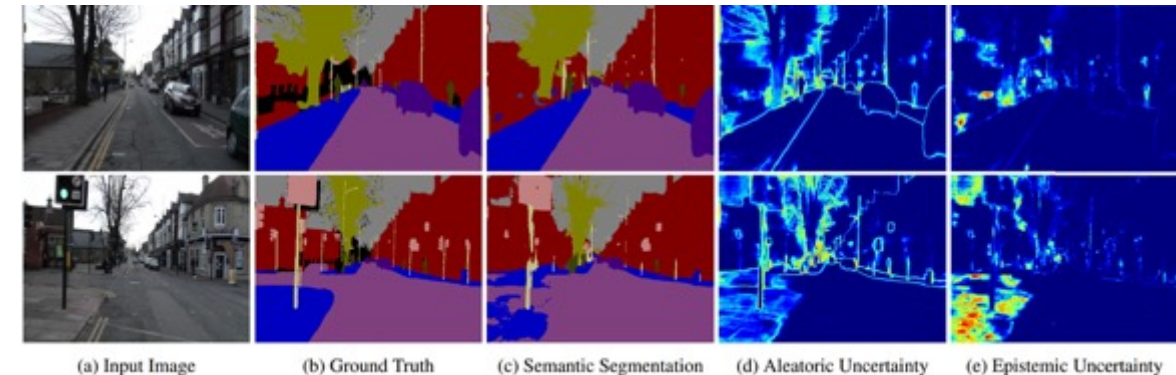
Surogate Model

Mark Craven, Jude W. Shavlik: Extracting Tree-Structured Representations of Trained Networks. NIPS 1995: 24-30



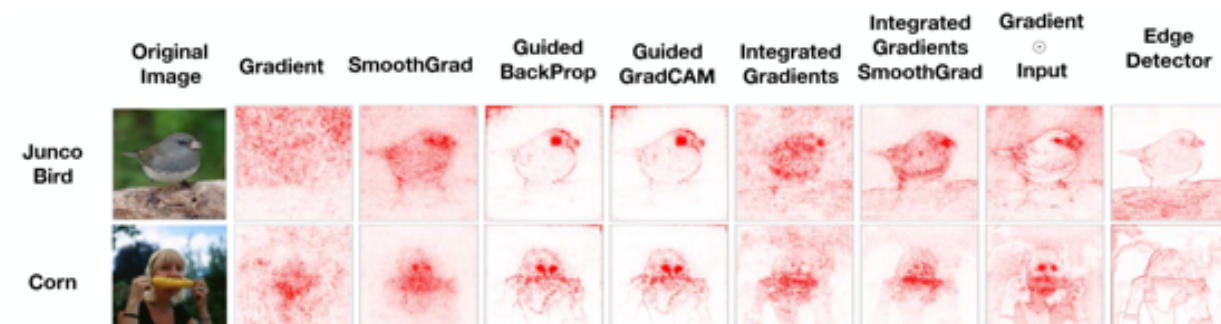
# Explanation in different AI fields

- Machine Learning
- Computer Vision



## Uncertainty Map

Alex Kendall, Yarin Gal: What Uncertainties Do We Need in Bayesian Deep Learning for Computer Vision? NIPS 2017: 5580-5590

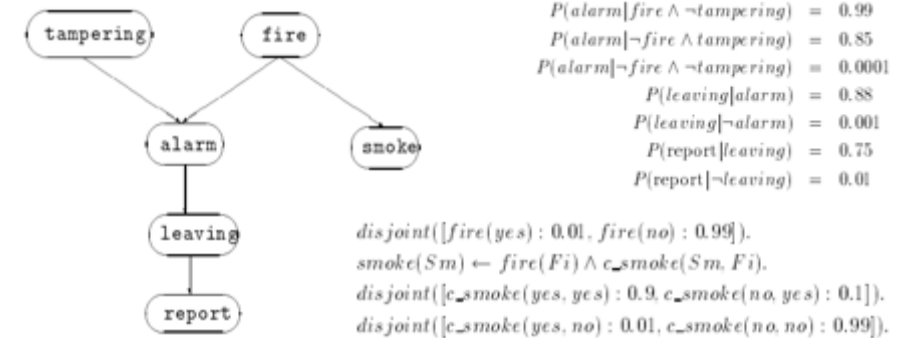


## Saliency Map

Julius Adebayo, Justin Gilmer, Michael Muelly, Ian J. Goodfellow, Moritz Hardt, Been Kim: Sanity Checks for Saliency Maps. NeurIPS 2018: 9525-9536

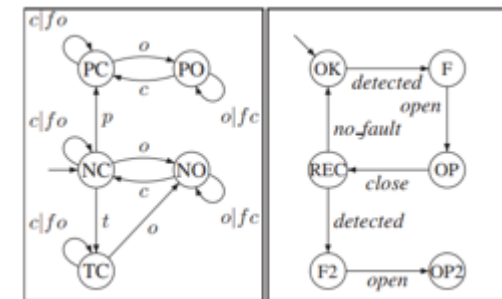
# Explanation in different AI fields

- Machine Learning
- Computer Vision
- Knowledge Representation and Reasoning



## Abduction Reasoning (in Bayesian Network)

David Poole: Probabilistic Horn Abduction and Bayesian Networks. *Artif. Intell.* 64(1): 81-129 (1993)

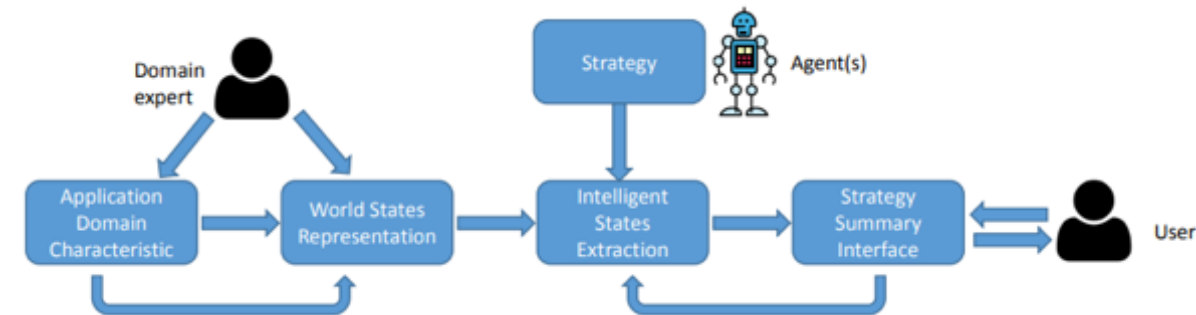


## Diagnosis Inference

Alban Grastien, Patrik Haslum, Sylvie Thiébaux: Conflict-Based Diagnosis of Discrete Event Systems: Theory and Practice. *KR* 2012

# Explanation in different AI fields

- Machine Learning
- Computer Vision
- Knowledge Representation and Reasoning
- Multi-agent Systems



## Agent Strategy Summarization

Ofra Amir, Finale Doshi-Velez, David Sarne: Agent Strategy Summarization. AAMAS 2018: 1203-1207

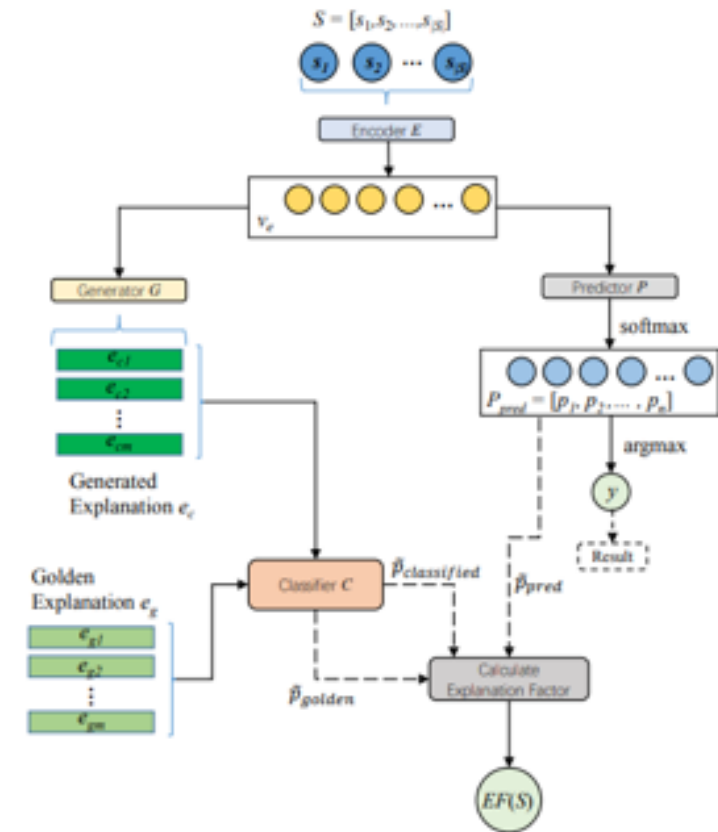


## Explainable Agents

Joost Broekens, Maaïke Harbers, Koen V. Hindriks, Karel van den Bosch, Catholijn M. Jonker, John-Jules Ch. Meyer: Do You Get It? User-Evaluated Explainable BDI Agents. MATES 2010: 28-39

# Explanation in different AI fields

- Machine Learning
- Computer Vision
- Knowledge Representation and Reasoning
- Multi-agent Systems
- NLP

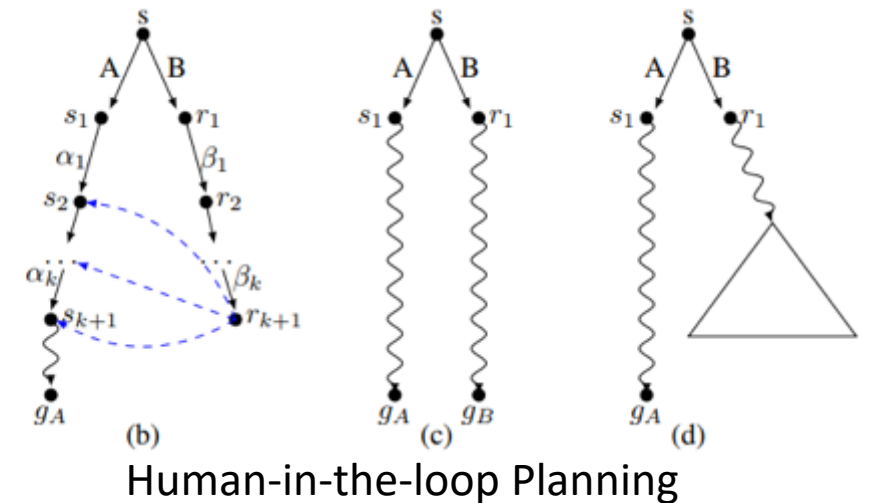


Explainable NLP

Hui Liu, Qingyu Yin, William Yang Wang: Towards Explainable NLP: A Generative Explanation Framework for Text Classification. CoRR abs/1811.00196 (2018)

# Explanation in different AI fields

- Machine Learning
- Computer Vision
- Knowledge Representation and Reasoning
- Multi-agent Systems
- NLP
- Planning and Scheduling



Maria Fox, Derek Long, Daniele Magazzeni: Explainable Planning. CoRR abs/1709.10256 (2017)

# Explanation in different AI fields

- Machine Learning
- Computer Vision
- Knowledge Representation and Reasoning
- Multi-agent Systems
- NLP
- Planning and Scheduling
- Robotics

**Robot:** I have decided to turn left.

**Human:** Why did you do that?

**Robot:** I believe that the correct action is to turn left  
BECAUSE:  
I'm being asked to go forward  
AND This area in front of me was 20 cm higher than me  
\*highlights area\*  
AND the area to the left has maximum protrusions of less than 5 cm \*highlights area\*  
AND I'm tilted to the right by more than 5 degrees.  
Here is a display of the path through the tree that lead to this decision. \*displays tree\*

**Human:** How confident are you in this decision?

**Robot:** The distribution of actions that reached this leaf node is shown in this histogram. \*displays histogram\*  
This action is predicted to be correct 67% of the time.

**Human:** Where did the threshold for the area in front come from?

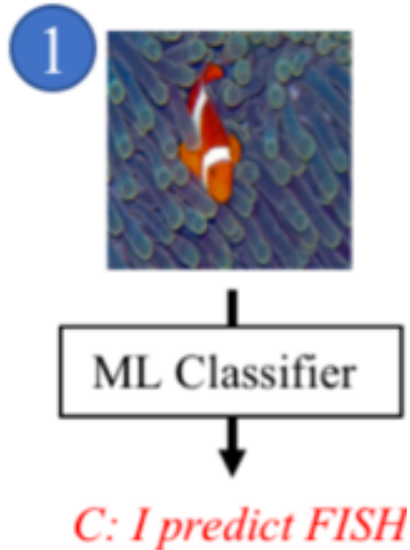
**Robot:** Here is the histogram of all training examples that reached this leaf. 80% of examples where this area was above 20 cm predicted the appropriate action to be "drive forward".

From Decision Tree to human-friendly information

Raymond Ka-Man Sheh: "Why Did You Do That?" Explainable Intelligent Robots. AAAI Workshops 2017

# Explanation as *Machine-Human Conversation*

[Weld and Bansal 2018]



- Humans may have follow-up questions
- Explanations cannot answer all users' concerns



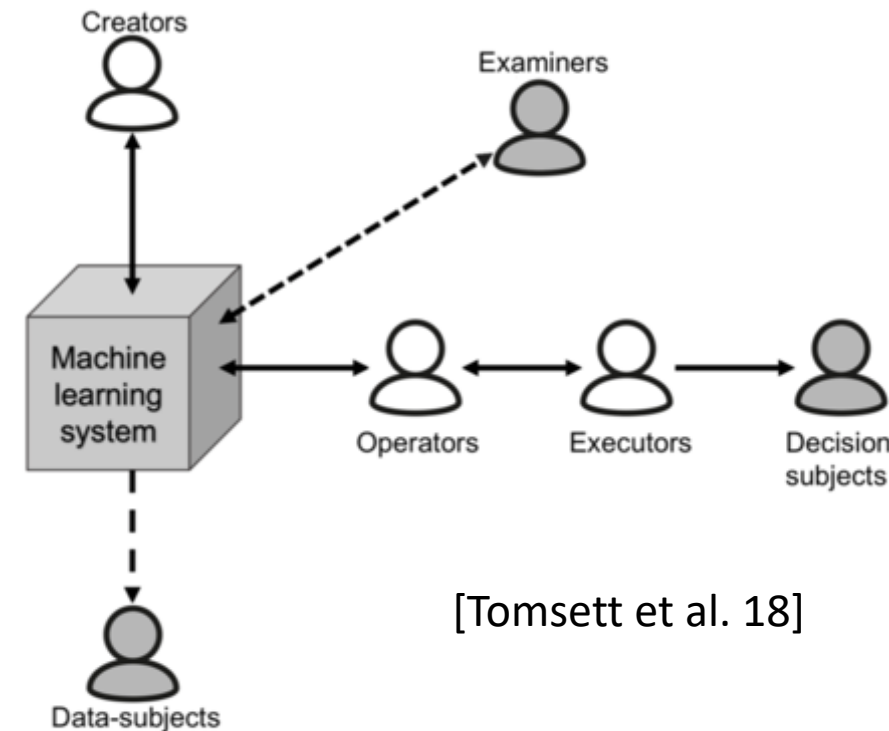
# Role-based Interpretability

~~“Is the explanation interpretable?”~~ → “*To whom* is the explanation interpretable?”

No Universally Interpretable Explanations!

- **End users** “Am I being treated fairly?”  
“Can I contest the decision?”  
“What could I do differently to get a positive outcome?”
- **Engineers, data scientists:** “Is my system working as designed?”
- **Regulators** “Is it compliant?”

An ideal explainer should model the *user background*.



[Tomsett et al. 18]

# Summarizing: the Need to Explain comes from ...

---

- User Acceptance & Trust

[Lipton 2016, Ribeiro 2016, Weld and Bansal 2018]

- Legal

- Conformance to ethical standards, fairness
- *Right to be informed*
- Contestable decisions

[Goodman and Flaxman 2016, Wachter 2017]

- Explanatory Debugging

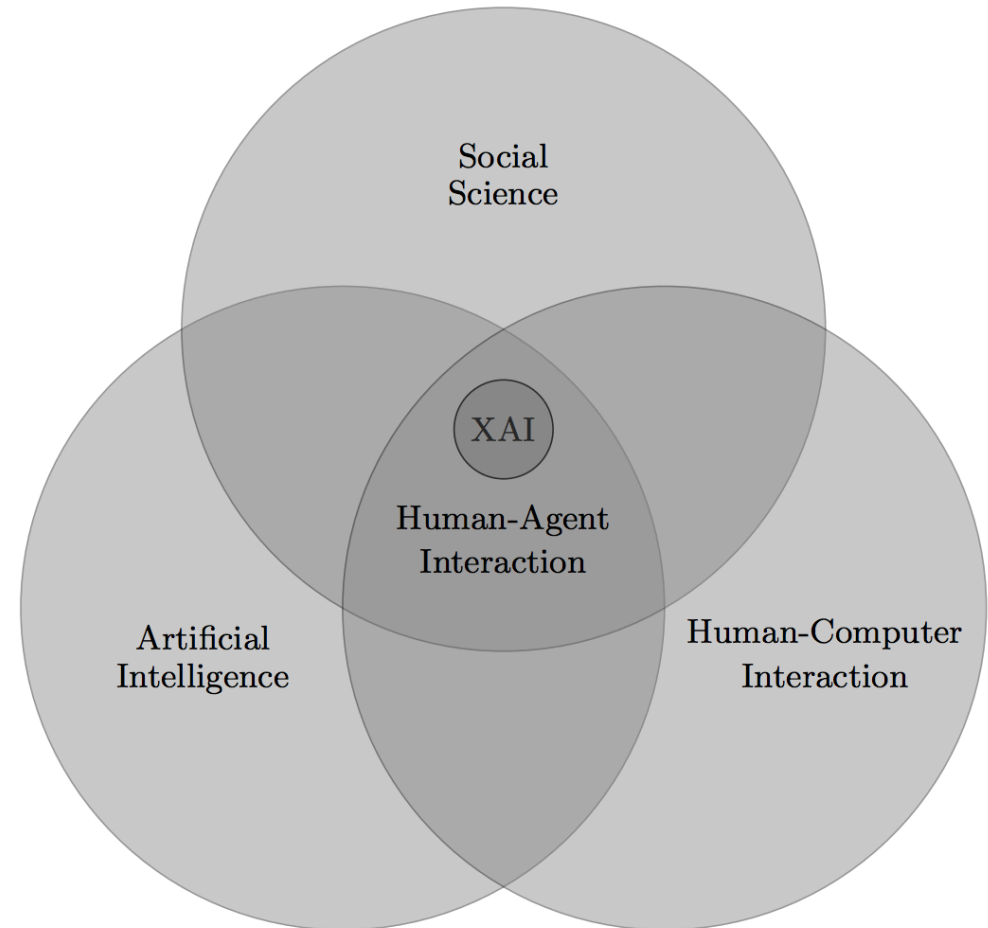
- Flawed performance metrics
- Inadequate features
- Distributional drift

[Kulesza et al. 2014, Weld and Bansal 2018]

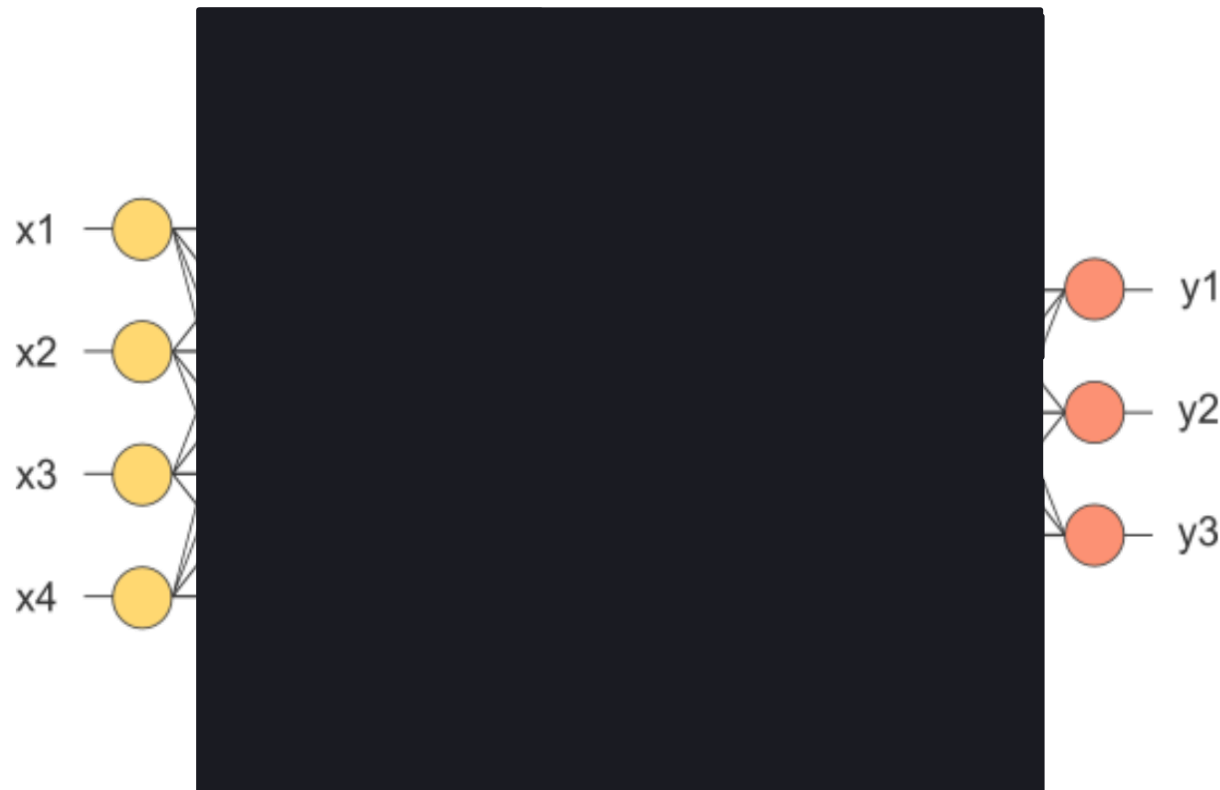
# XAI is Interdisciplinary

---

- For millennia, philosophers have asked the questions about what constitutes an explanation, what is the function of explanations, and what are their structure
- **[Tim Miller 2018]**



# What is a Black Box Model?



A ***black box*** is a model, whose internals are either unknown to the observer or they are known but uninterpretable by humans.

- Guidotti, R., Monreale, A., Ruggieri, S., Turini, F., Giannotti, F., & Pedreschi, D. (2018). *A survey of methods for explaining black box models*. *ACM Computing Surveys (CSUR)*, 51(5), 93.






Needs For Interpretable Models



# COMPAS recidivism black bias



DYLAN FUGETT

Prior Offense  
1 attempted burglary

Subsequent Offenses  
3 drug possessions

LOW RISK

3



BERNARD PARKER


Prior Offense  
1 resisting arrest  
without violence

Subsequent Offenses  
None

HIGH RISK

10

*Fugett was rated low risk after being arrested with cocaine and marijuana. He was arrested three times on drug charges after that.*





H

H

W

W

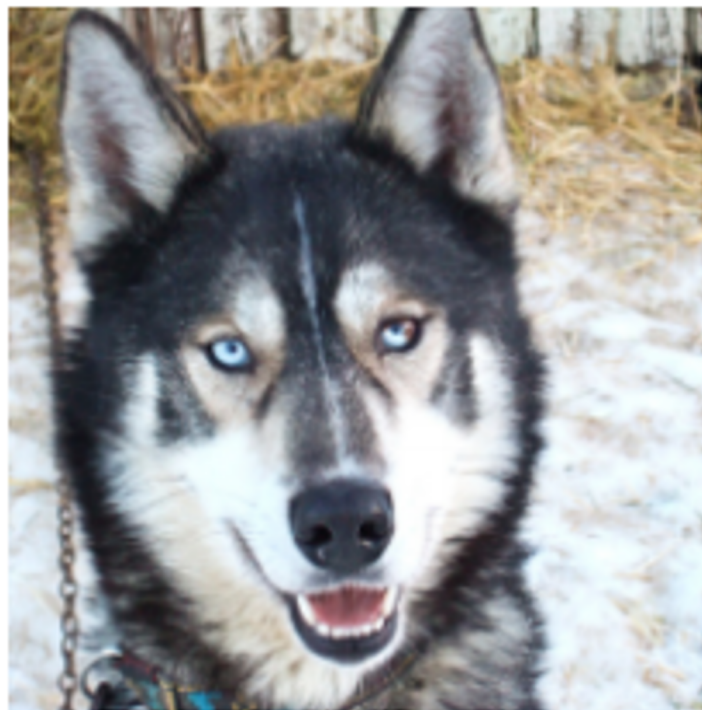
# The background bias



H



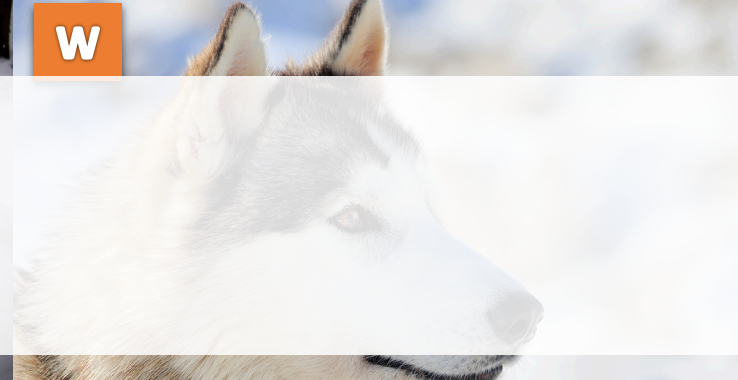
H



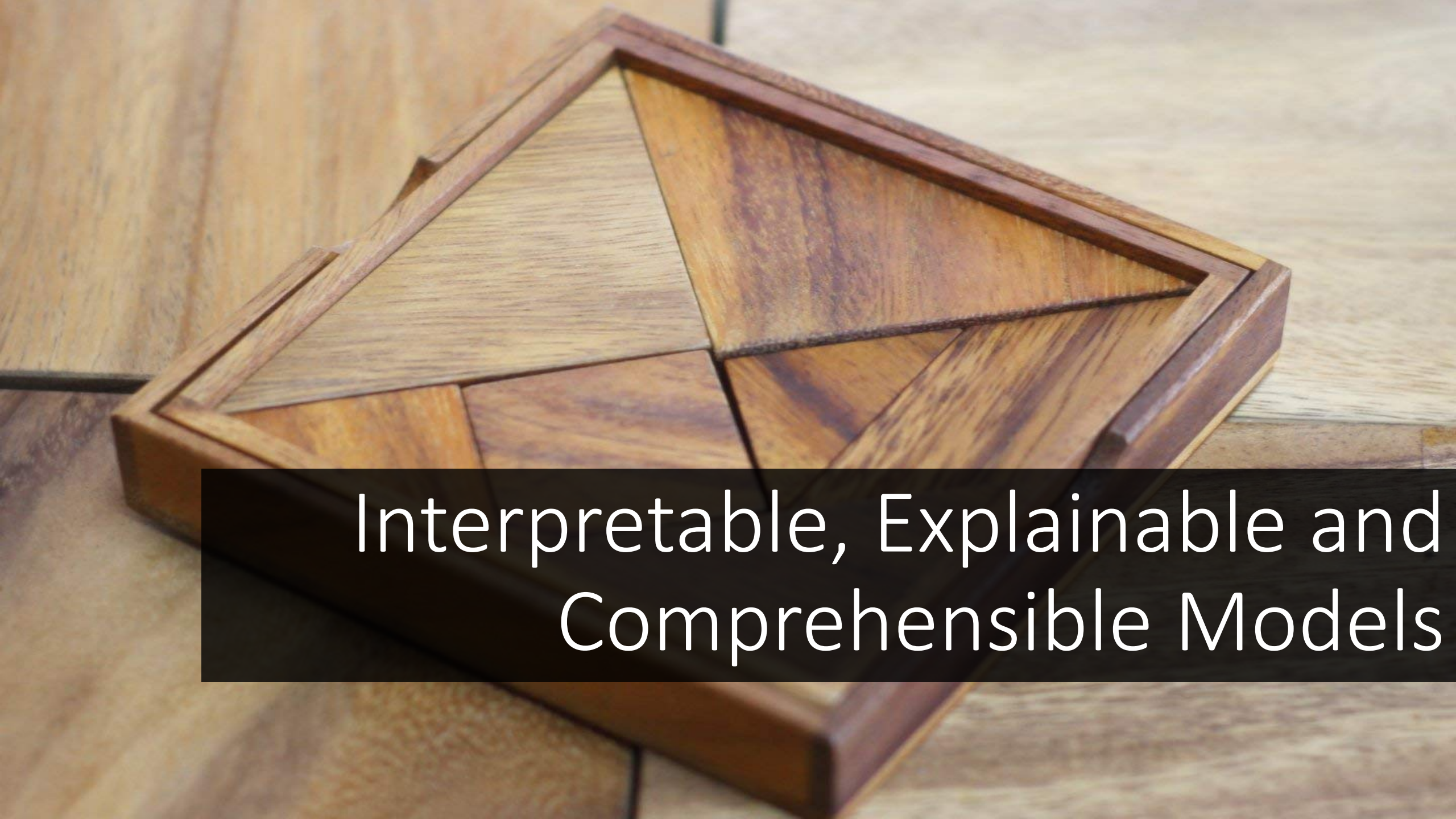
(a) Husky classified as wolf



(b) Explanation







Interpretable, Explainable and  
Comprehensible Models



# Interpretability

---

- To ***interpret*** means to give or provide the meaning or to explain and present in understandable terms some concepts.
- In data mining and machine learning, interpretability is the ***ability to explain*** or to provide the meaning ***in understandable terms to a human***.

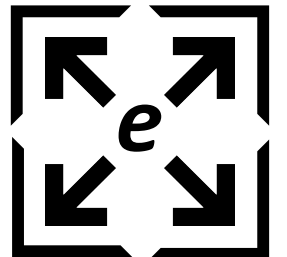


- <https://www.merriam-webster.com/>
- Finale Doshi-Velez and Been Kim. 2017. ***Towards a rigorous science of interpretable machine learning***. arXiv:1702.08608v2.

# Dimensions of Interpretability

---

- ***Global and Local Interpretability:***
  - *Global*: understanding the whole logic of a model
  - *Local*: understanding only the reasons for a specific decision
- ***Time Limitation:*** the time that the user can spend for understanding an explanation.
- ***Nature of User Expertise:*** users of a predictive model may have different background knowledge and experience in the task. The nature of the user expertise is a key aspect for interpretability of a model.



# Desiderata of an Interpretable Model

---

- ***Interpretability*** (or comprehensibility): to which extent the model and/or its predictions are human understandable. Is measured with the *complexity* of the model.
- ***Fidelity***: to which extent the model imitate a black-box predictor.
- ***Accuracy***: to which extent the model predicts unseen instances.

- Alex A. Freitas. 2014. ***Comprehensible classification models: A position paper***. ACM SIGKDD Explor. Newslett.



# Desiderata of an Interpretable Model

---

- ***Fairness***: the model guarantees the protection of groups against discrimination.
- ***Privacy***: the model does not reveal sensitive information about people.
- ***Respect Monotonicity***: the increase of the values of an attribute either increase or decrease in a monotonic way the probability of a record of being member of a class.
- ***Usability***: an interactive and queryable explanation is more usable than a textual and fixed explanation.

- Andrea Romei and Salvatore Ruggieri. 2014. ***A multidisciplinary survey on discrimination analysis***. Knowl. Eng.
- Yousra Abdul Alsahib S. Aldeen, Mazleena Salleh, and Mohammad Abdur Razzaque. 2015. ***A comprehensive review on privacy preserving data mining***. SpringerPlus .
- Alex A. Freitas. 2014. ***Comprehensible classification models: A position paper***. ACM SIGKDD Explor. Newslett.



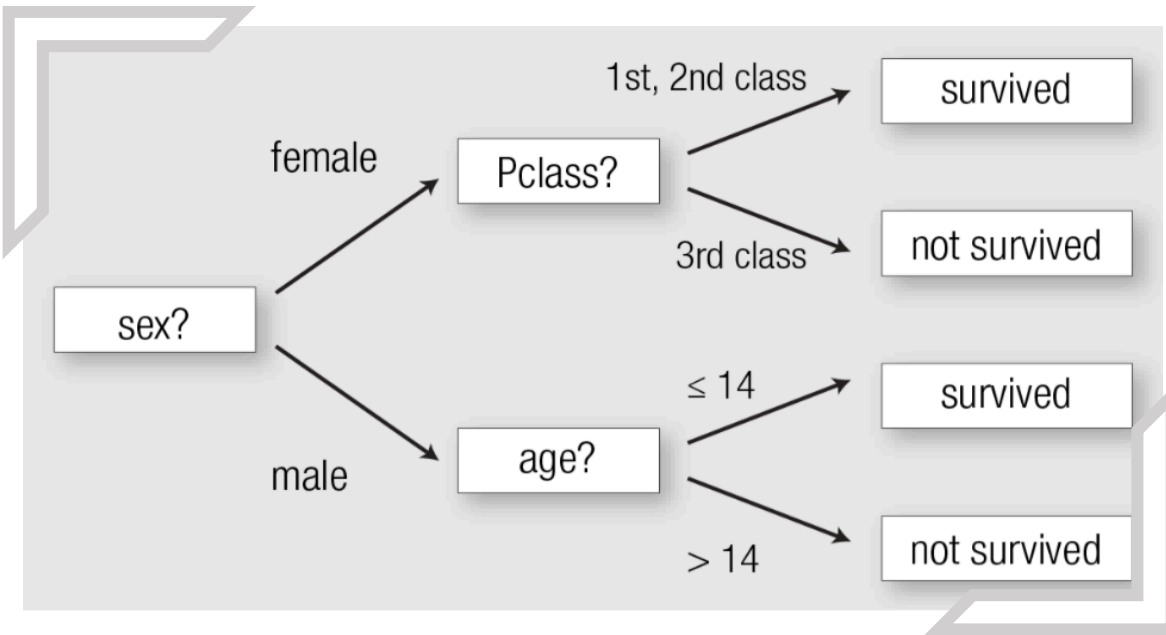
# Desiderata of an Interpretable Model

---

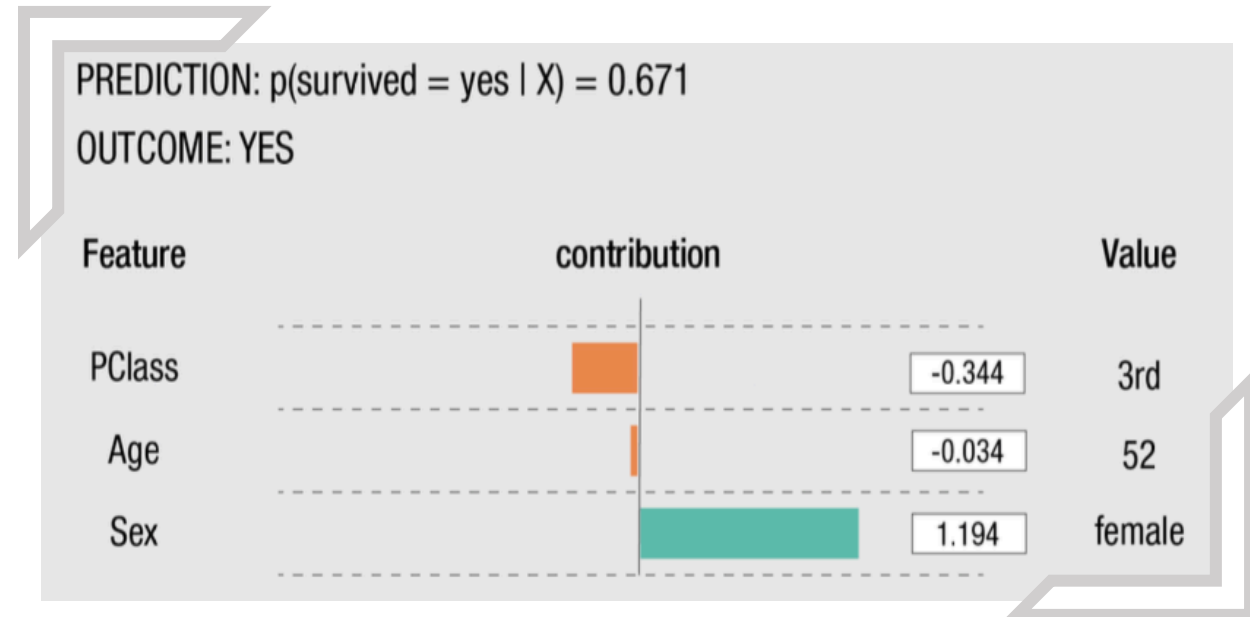
- ***Reliability and Robustness***: the interpretable model should maintain high levels of performance independently from small variations of the parameters or of the input data.
- ***Causality***: controlled changes in the input due to a perturbation should affect the model behavior.
- ***Scalability***: the interpretable model should be able to scale to large input data with large input spaces.
- ***Generality***: the model should not require special training or restrictions.



# Recognized Interpretable Models



Decision Tree

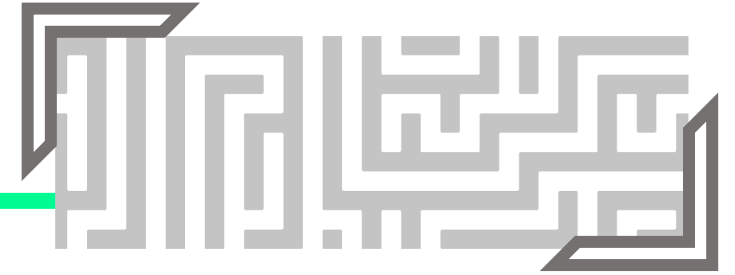


Linear Model

*if condition<sub>1</sub>  $\wedge$  condition<sub>2</sub>  $\wedge$  condition<sub>3</sub> then outcome*

Rules

# Complexity



- Opposed to *interpretability*.
- Is only related to the model and not to the training data that is unknown.
- Generally estimated with a rough approximation related to the **size** of the interpretable model.
- Linear Model: number of non zero weights in the model.
- Rule: number of attribute-value pairs in condition.
- Decision Tree: estimating the complexity of a tree can be hard.

- Marco Tulio Ribeiro, Sameer Singh, and Carlos Guestrin. 2016. ***Why should i trust you?: Explaining the predictions of any classifier***. KDD.
- Houtao Deng. 2014. ***Interpreting tree ensembles with intrees***. arXiv preprint arXiv:1408.5456.
- Alex A. Freitas. 2014. ***Comprehensible classification models: A position paper***. ACM SIGKDD Explor. Newslett.

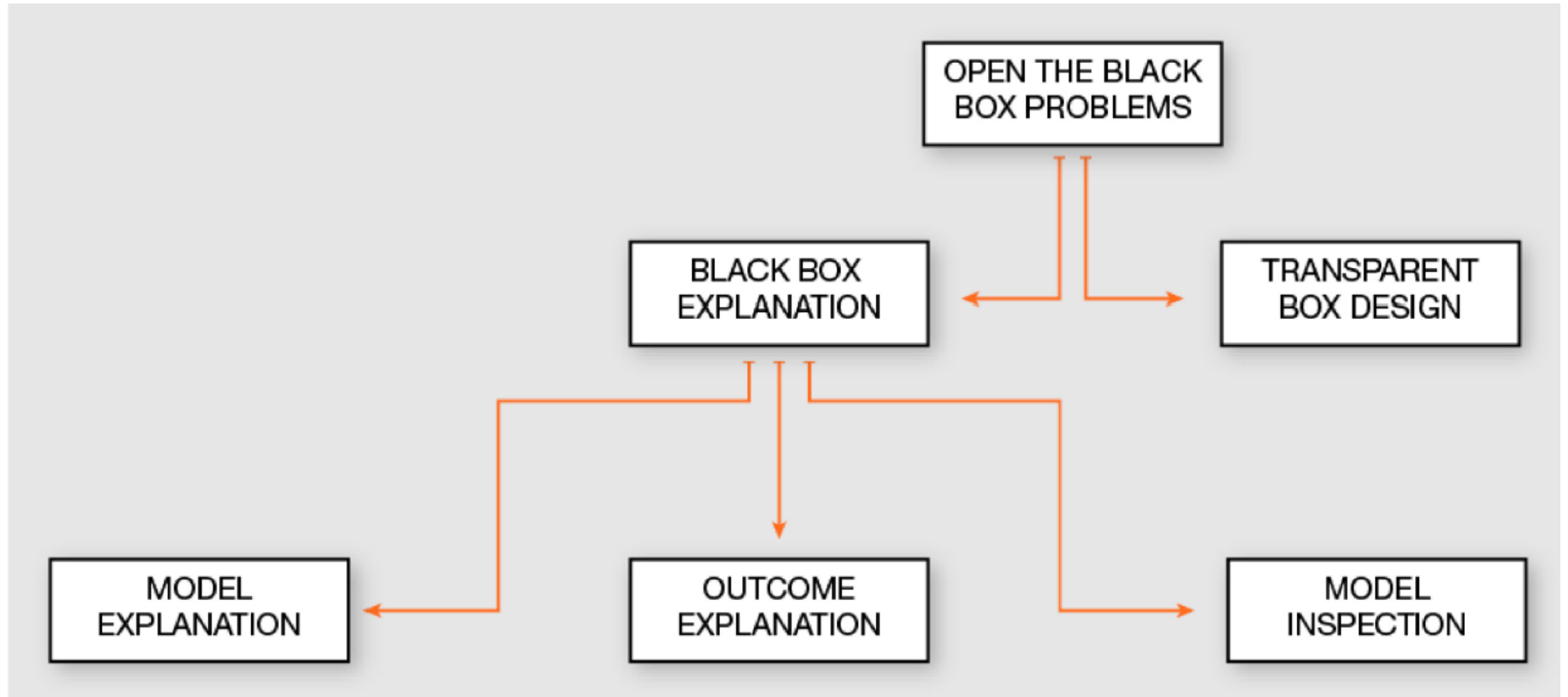


A close-up photograph of a person's hand turning a silver-colored dial on a dark, textured surface, likely a safe or a black box. The dial has numbers 60, 70, 80, and 90 visible. A key is held in the bottom left corner. A semi-transparent black banner with white text is overlaid across the bottom half of the image.

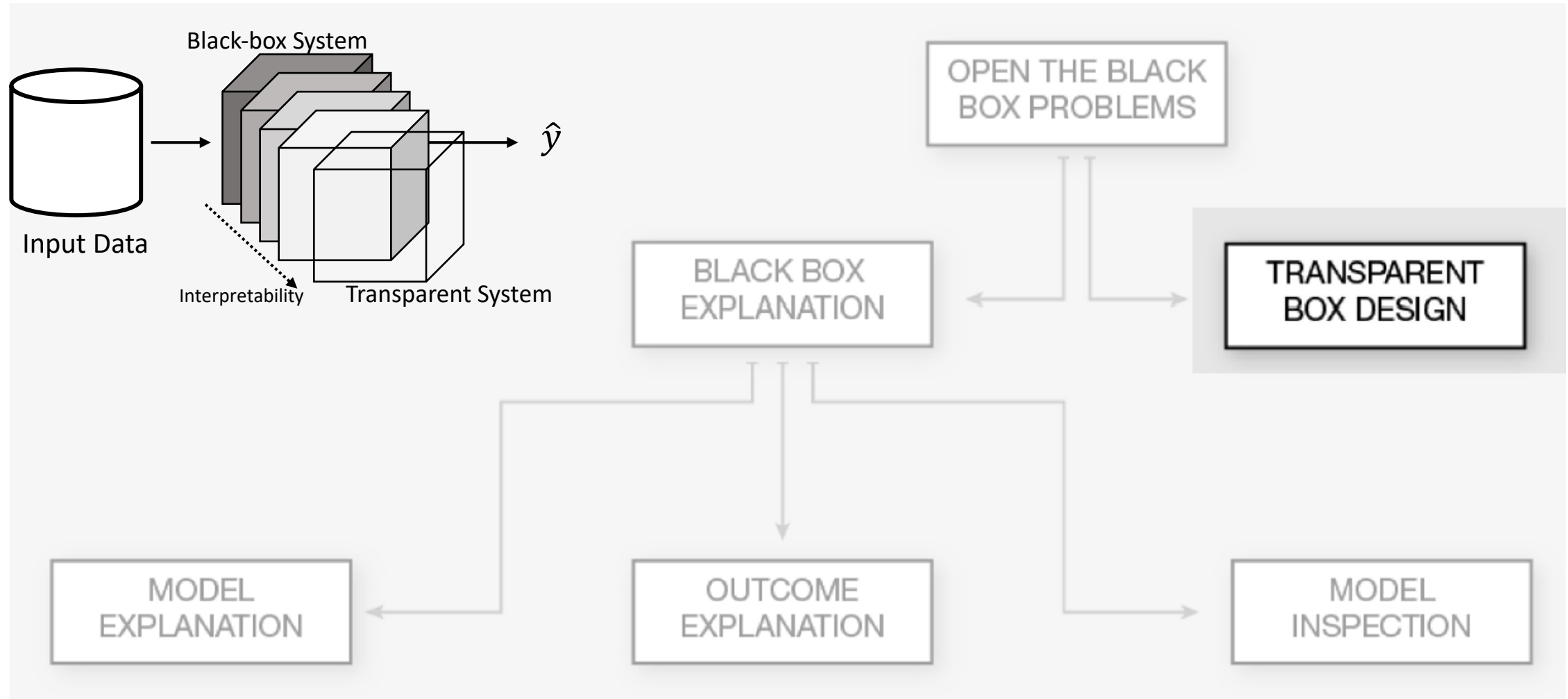
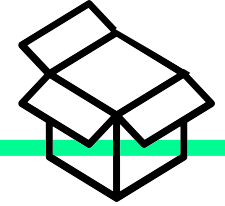
Open the Black Box Problems



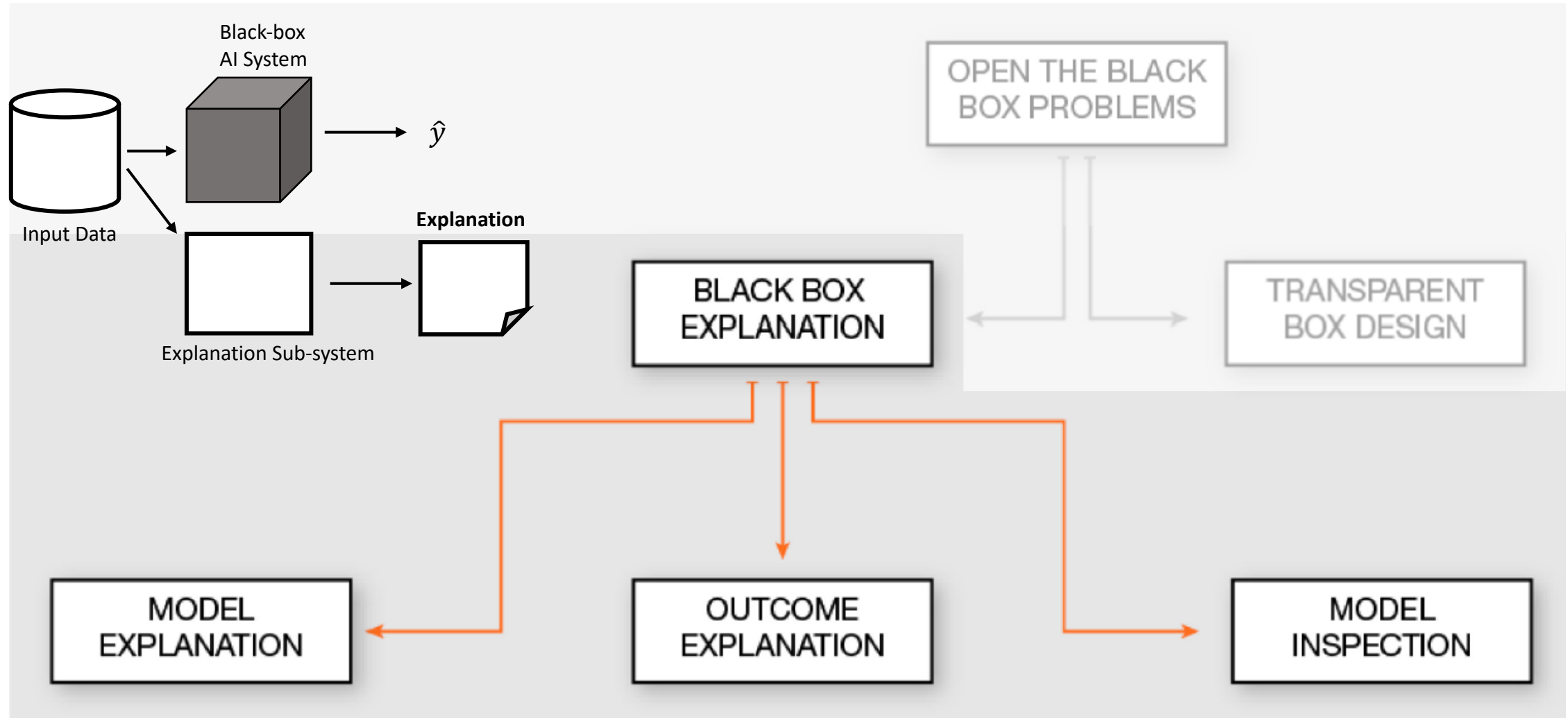
# Problems Taxonomy



# XbD – eXplanation by Design

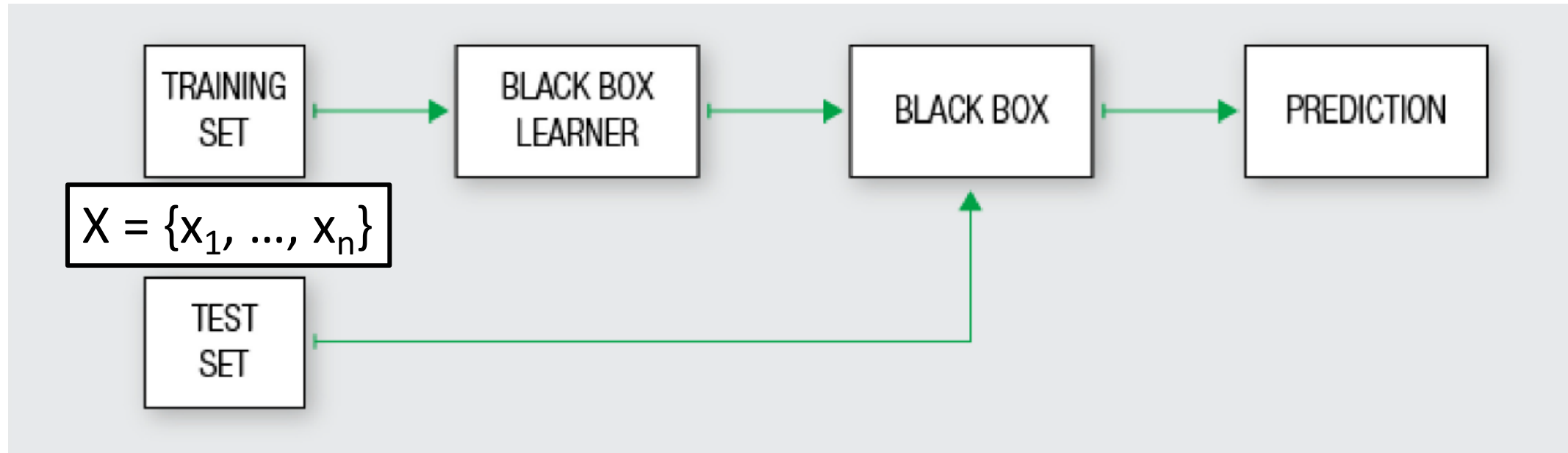


# BBX - Black Box eXplanation



# Classification Problem

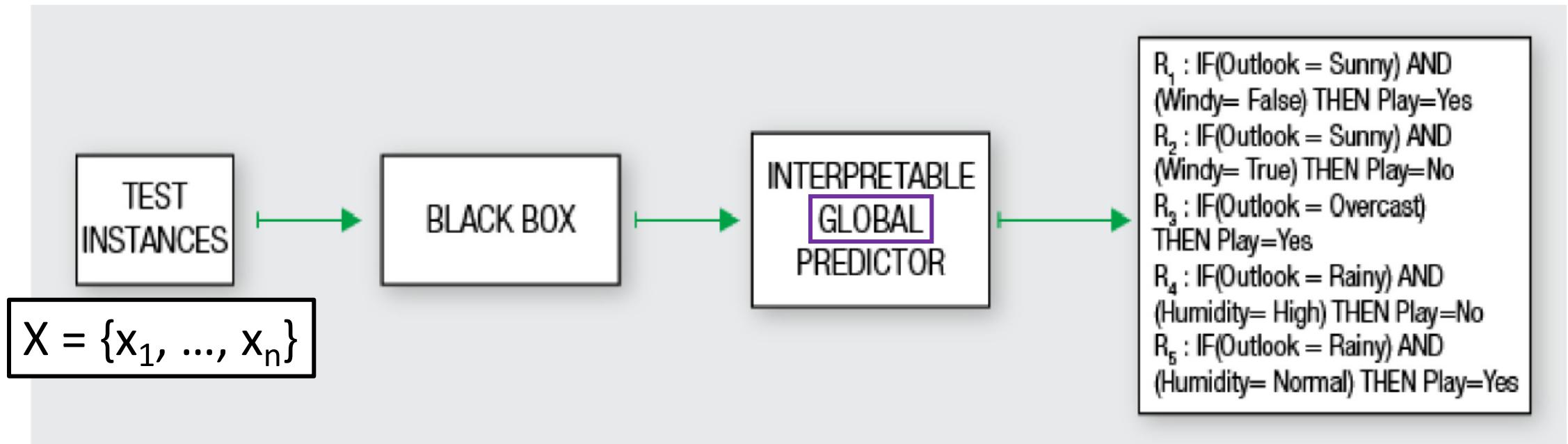
---



# Model Explanation Problem



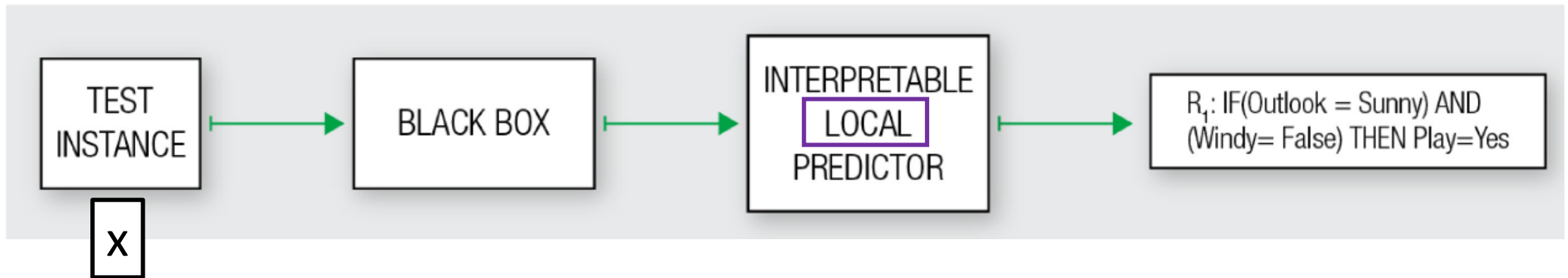
Provide an interpretable model able to mimic the ***overall logic/behavior*** of the black box and to explain its logic.



# Outcome Explanation Problem



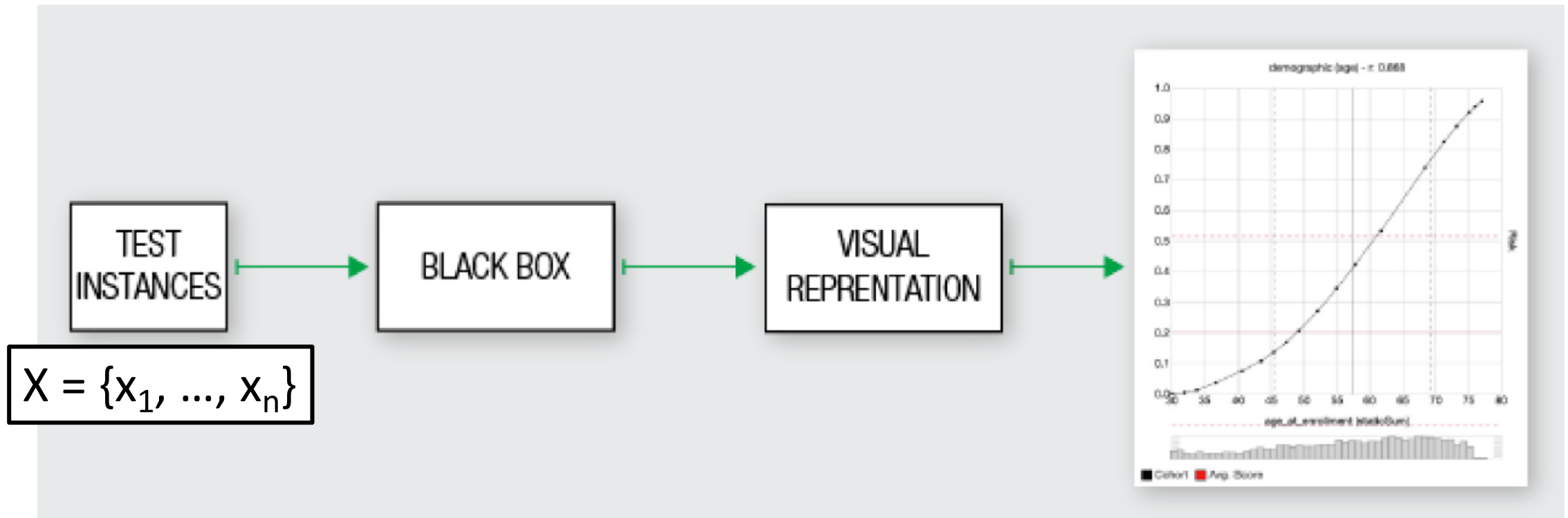
Provide an interpretable outcome, i.e., an ***explanation*** for the outcome of the black box for a ***single instance***.



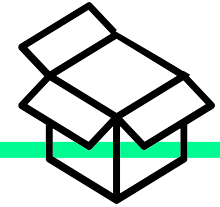
# Model Inspection Problem



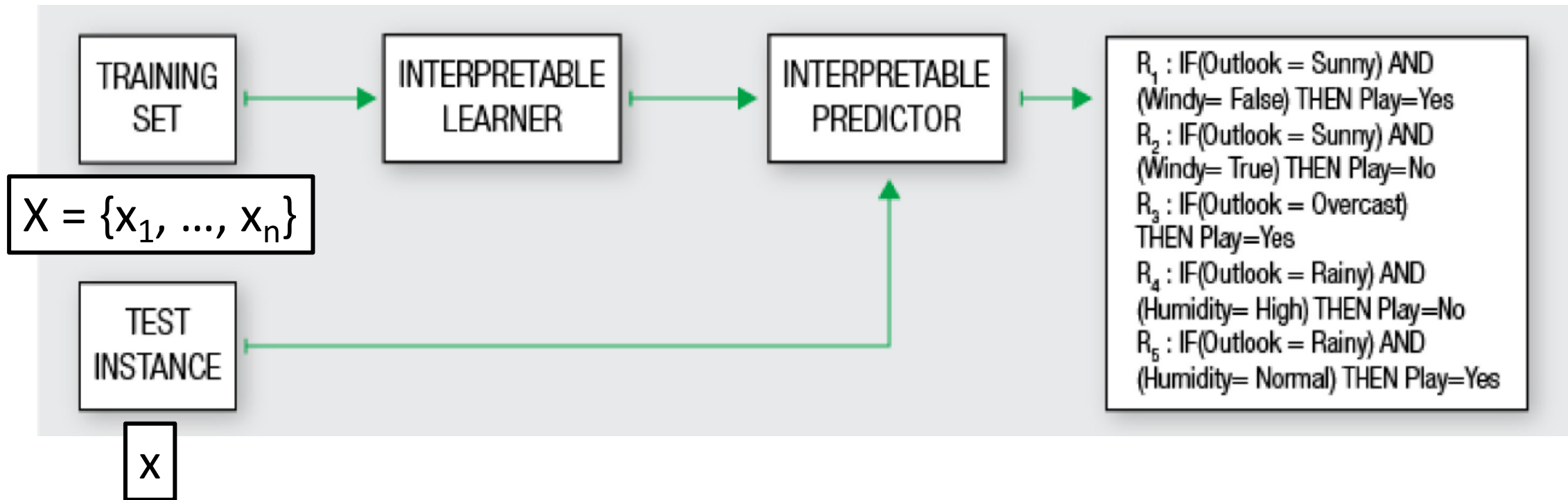
Provide a representation (visual or textual) for understanding either how the black box model works or why the black box returns certain predictions more likely than others.



# Transparent Box Design Problem



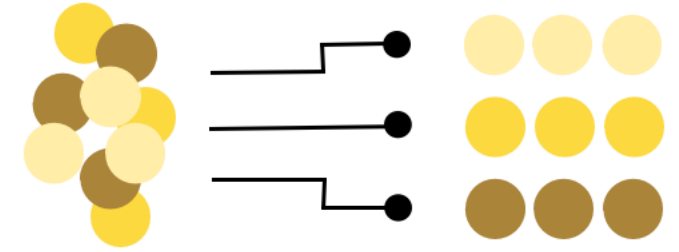
Provide a model which is locally or globally interpretable on its own.





# Categorization

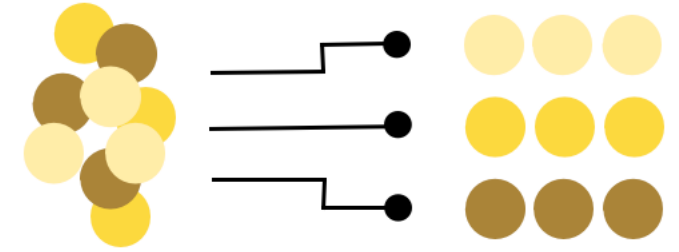
---



- The type of ***problem***
- The type of ***black box model*** that the explainer is able to open
- The type of ***data*** used as input by the black box model
- The type of ***explainer*** adopted to open the black box

# Black Boxes

---



- Neural Network (***NN***)
- Tree Ensemble (***TE***)
- Support Vector Machine (***SVM***)
- Deep Neural Network (***DNN***)



name	rank	gender	year
Jacob	1	boy	2010
Isabella	1	girl	2010
Ethan	2	boy	2010
Sophia	2	girl	2010
Michael	3	boy	2010

One row  
(4 fields)

# Tabular (TAB)

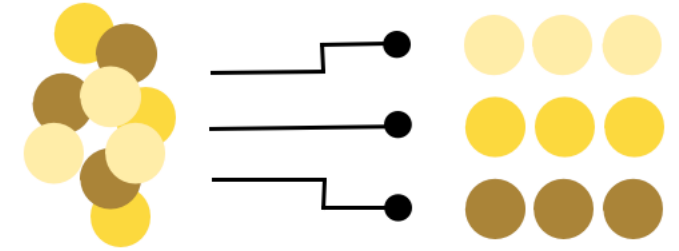
## A collage of three photographs. The top-left photo shows a young girl with blonde hair, wearing a yellow patterned dress and blue leggings, riding a blue tricycle on a paved path. The top-right photo shows two white Arctic wolves standing on a rocky, grassy slope. The bottom photo shows a reddish-brown fox walking in a dry, grassy field. The text '© Bernard Castelein / naturepl.com' is visible in the bottom right corner of the fox photo.

## A large pile of small, rectangular word cards scattered on a white background. The cards are white with black text. The words are in various orientations, some upright and some upside down. Visible words include: 'picture', 'behind', 'your', 'spring', 'time', 'size', 'whisper', 'chime', 'sing', 'moment', 'leave', 'purple', 'delicious', 'visit', 'taste', 'would', 'the', 'world', 'may', 'beauty', 'for', 'in', 'like', 'an', 'pc', 'behind', 'your', 'picture'. The cards are piled together, creating a sense of abundance and randomness.

# Explainers

---

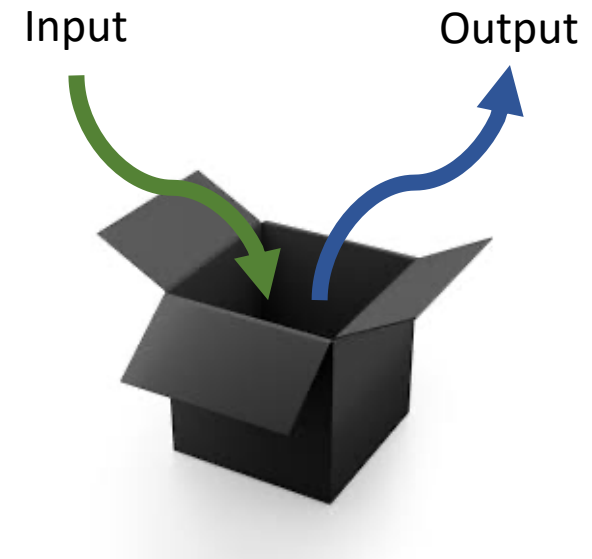
- Decision Tree (**DT**)
- Decision Rules (**DR**)
- Features Importance (**FI**)
- Saliency Maps (**SM**)
- Sensitivity Analysis (**SA**)
- Partial Dependence Plot (**PDP**)
- Prototype Selection (**PS**)
- Activation Maximization (**AM**)



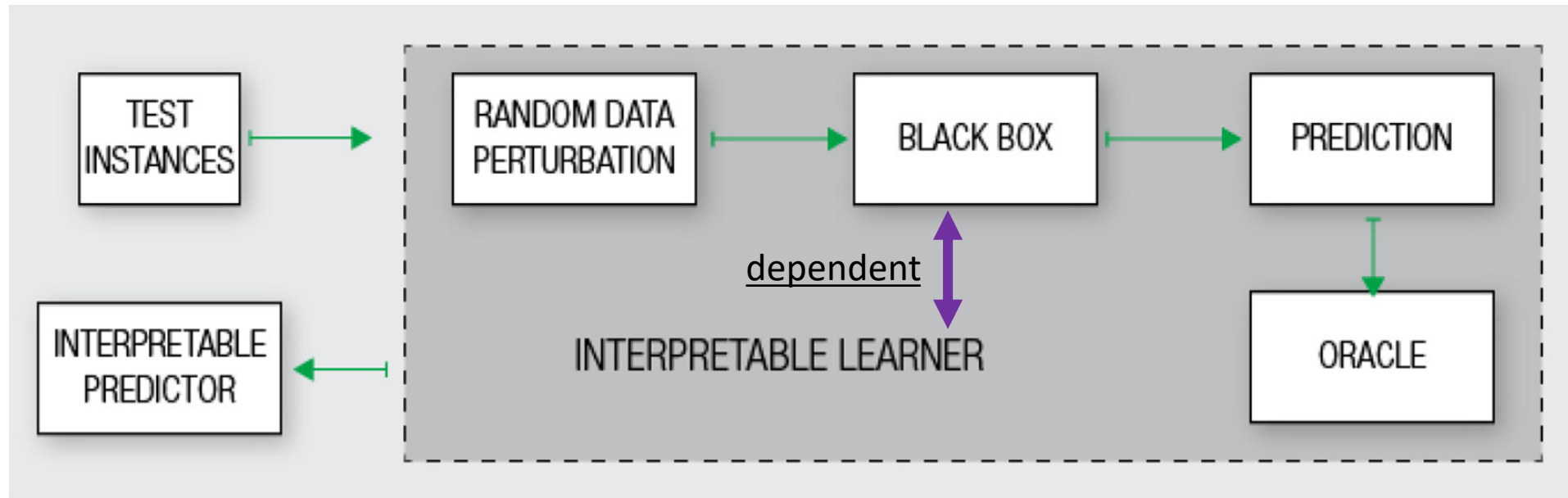
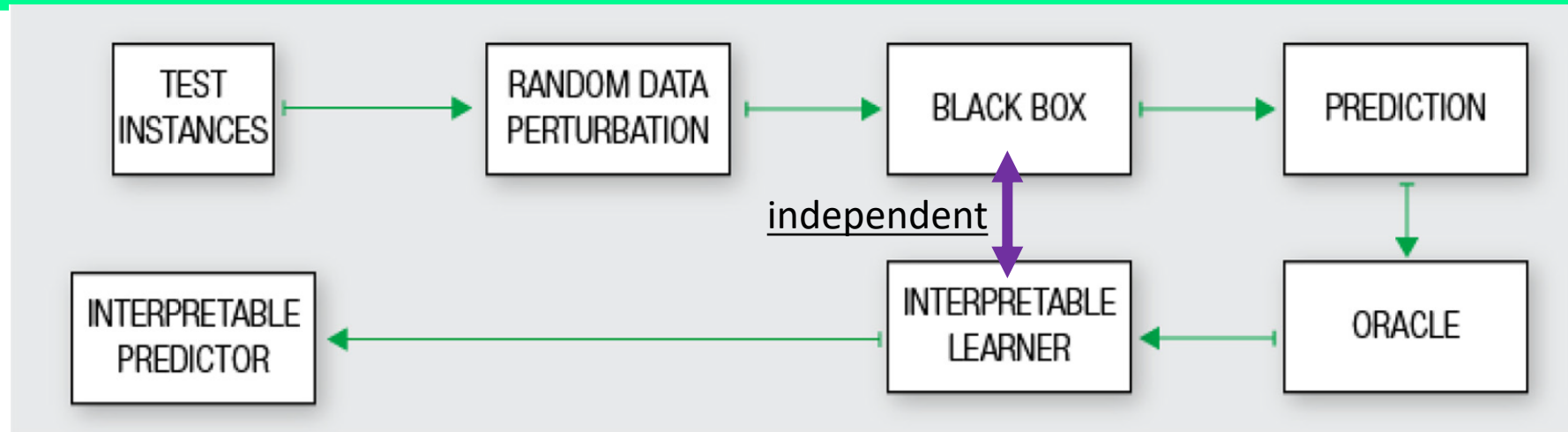
# Reverse Engineering

---

- The name comes from the fact that we can only **observe** the **input** and **output** of the black box.
- Possible actions are:
  - **choice** of a particular comprehensible predictor
  - querying/auditing the black box with input records created in a controlled way using **random perturbations** w.r.t. a certain prior knowledge (e.g. train or test)
- It can be **generalizable or not**:
  - Model-Agnostic
  - Model-Specific



# Model-Agnostic vs Model-Specific



<i>Name</i>	<i>Ref.</i>	<i>Authors</i>	<i>Year</i>	<i>Explanator</i>	<i>Black Box</i>	<i>Data Type</i>	<i>General</i>	<i>Random</i>	<i>Examples</i>	<i>Code</i>	<i>Dataset</i>
Trepan	[22]	Craven et al.	1996	DT	NN	TAB	✓				✓
—	[57]	Krishnan et al.	1999	DT	NN	TAB	✓		✓		✓
DecText	[12]	Boz	2002	DT	NN	TAB	✓	✓			✓
GPDT	[46]	Johansson et al.	2009	DT	NN	TAB	✓	✓	✓		✓
Tree Metrics	[17]	Chipman et al.	1998	DT	TE	TAB					✓
CCM	[26]	Domingos et al.	1998	DT	TE	TAB	✓	✓			✓
—	[34]	Gibbons et al.	2013	DT	TE	TAB	✓	✓			
STA	[140]	Zhou et al.	2016	DT	TE	TAB		✓			
CDT	[104]	Schetinin et al.	2007	DT	TE	TAB			✓		
—	[38]	Hara et al.	2016	DT	TE	TAB		✓	✓		✓
TSP	[117]	Tan et al.	2016	DT	TE	TAB					✓
Conj Rules	[21]	Craven et al.	1999	DR	NN	TAB					
G-REX	[44]	Johansson et al.	2003	DR	NN	TAB	✓	✓	✓		
REFNE	[141]	Zhou et al.	2003	DR	NN	TAB	✓	✓	✓		✓
RxREN	[6]	Augusta et al.	2012	DR	NN	TAB		✓	✓		✓

Solving The Model Explanation Problem

# Global Model Explainers

- Explinator: DT
  - Black Box: NN, TE
  - Data Type: TAB
- Explinator: DR
  - Black Box: NN, SVM, TE
  - Data Type: TAB
- Explinator: FI
  - Black Box: AGN
  - Data Type: TAB

```
R1 : IF(Outlook = Sunny) AND  
(Windy= False) THEN Play=Yes  
R2 : IF(Outlook = Sunny) AND  
(Windy= True) THEN Play=No  
R3 : IF(Outlook = Overcast)  
THEN Play=Yes  
R4 : IF(Outlook = Rainy) AND  
(Humidity= High) THEN Play=No  
R5 : IF(Outlook = Rainy) AND  
(Humidity= Normal) THEN Play=Yes
```

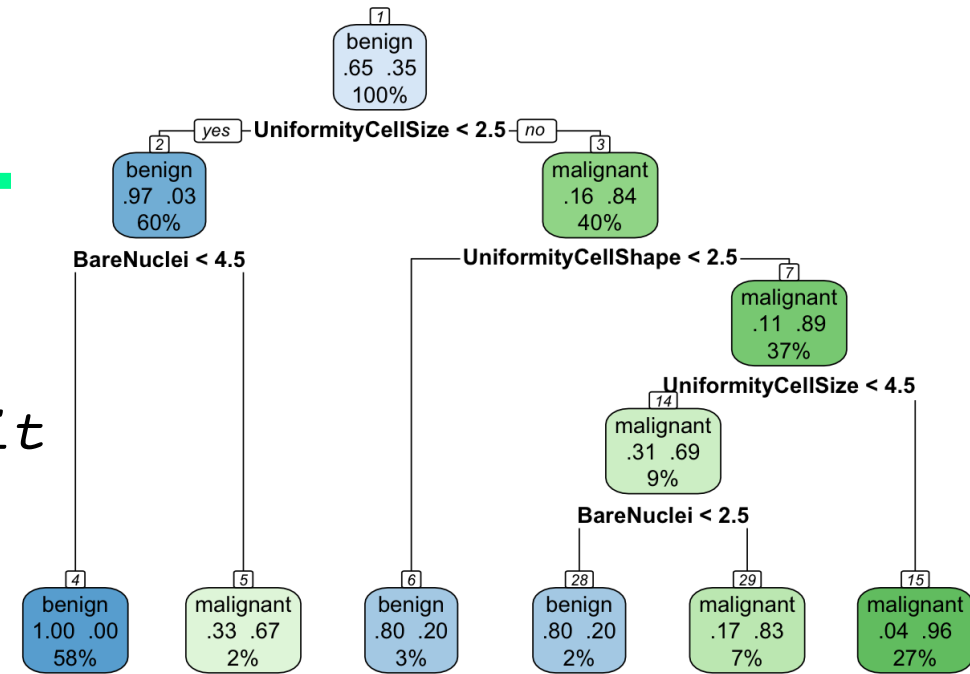


# Trepan – DT, NN, TAB

```

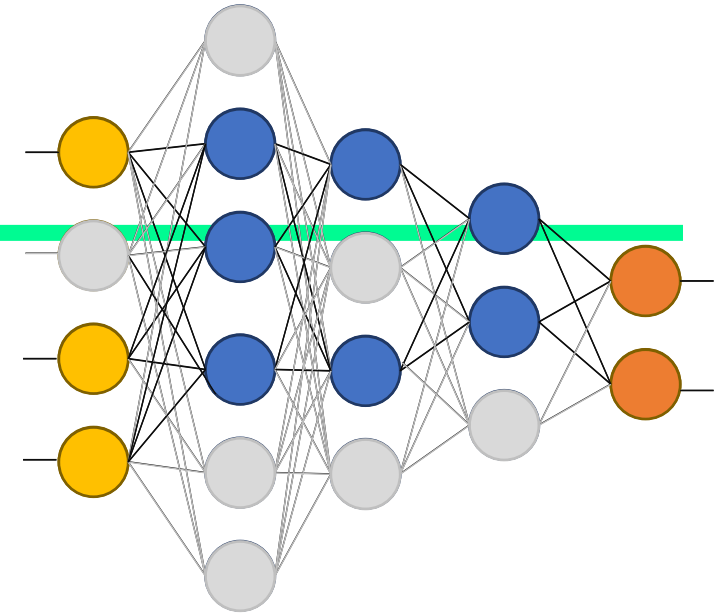
01  T = root_of_the_tree()
02  Q = <T, X, {}>
03  while Q not empty & size(T) < limit
04      N, XN, CN = pop(Q)
05      ZN = random(XN, CN)
06  black box auditing → yZ = b(Z), y = b(XN)
07      if same_class(y ∪ yZ)
08          continue
09      S = best_split(XN ∪ ZN, y ∪ yZ)
10      S' = best_m-of-n_split(S)
11      N = update_with_split(N, S')
12      for each condition c in S'
13          C = new_child_of(N)
14          CC = CN ∪ {c}
15          XC = select_with_constraints(XN, CN)
16          put(Q, <C, XC, CC>)

```



# RxREN – DR, NN, TAB

```
01  prune insignificant neurons
02  for each significant neuron
03    for each outcome
04    black box → compute mandatory data ranges
05    auditing
06    for each outcome
07      build rules using data ranges of each neuron
08  prune insignificant rules
09  update data ranges in rule conditions analyzing error
```



```
if ((data(I1) ≥ L13 ∧ data(I1) ≤ U13) ∧ (data(I2) ≥ L23 ∧ data(I2) ≤ U23) ∧
(data(I3) ≥ L33 ∧ data(I3) ≤ U33)) then class = C3
else
if ((data(I1) ≥ L11 ∧ data(I1) ≤ U11) ∧ (data(I3) ≥ L31 ∧ data(I3) ≤ U31))
then class = C1
else
class = C2
```

- M. Gethsiyal Augasta and T. Kathirvalavakumar. 2012.  
*Reverse engineering the neural networks for rule  
extraction in classification problems*. NPL.

<i>Name</i>	<i>Ref.</i>	<i>Authors</i>	<i>Year</i>	<i>Explanator</i>	<i>Black Box</i>	<i>Data Type</i>	<i>General</i>	<i>Random</i>	<i>Examples</i>	<i>Code</i>	<i>Dataset</i>
—	[134]	Xu et al.	2015	SM	DNN	IMG			✓	✓	✓
—	[30]	Fong et al.	2017	SM	DNN	IMG			✓		
CAM	[139]	Zhou et al.	2016	SM	DNN	IMG			✓	✓	✓
Grad-CAM	[106]	Selvaraju et al.	2016	SM	DNN	IMG			✓	✓	✓
—	[109]	Simonian et al.	2013	SM	DNN	IMG			✓		✓
PWD	[7]	Bach et al.	2015	SM	DNN	IMG			✓		✓
—	[113]	Sturm et al.	2016	SM	DNN	IMG			✓		✓
DTD	[78]	Montavon et al.	2017	SM	DNN	IMG			✓		✓
DeapLIFT	[107]	Shrikumar et al.	2017	FI	DNN	ANY			✓	✓	
CP	[64]	Landecker et al.	2013	SM	NN	IMG			✓		
—	[143]	Zintgraf et al.	2017	SM	DNN	IMG			✓	✓	✓
VBP	[11]	Bojarski et al.	2016	SM	DNN	IMG			✓		✓
—	[65]	Lei et al.	2016	SM	DNN	TXT			✓		✓
ExplainD	[89]	Poulin et al.	2006	FI	SVM	TAB		✓	✓		
—	[29]	Strumbelj et al.	2010	FI	AGN	TAB	✓	✓	✓		✓

# Solving The Outcome Explanation Problem

# Local Model Explainers

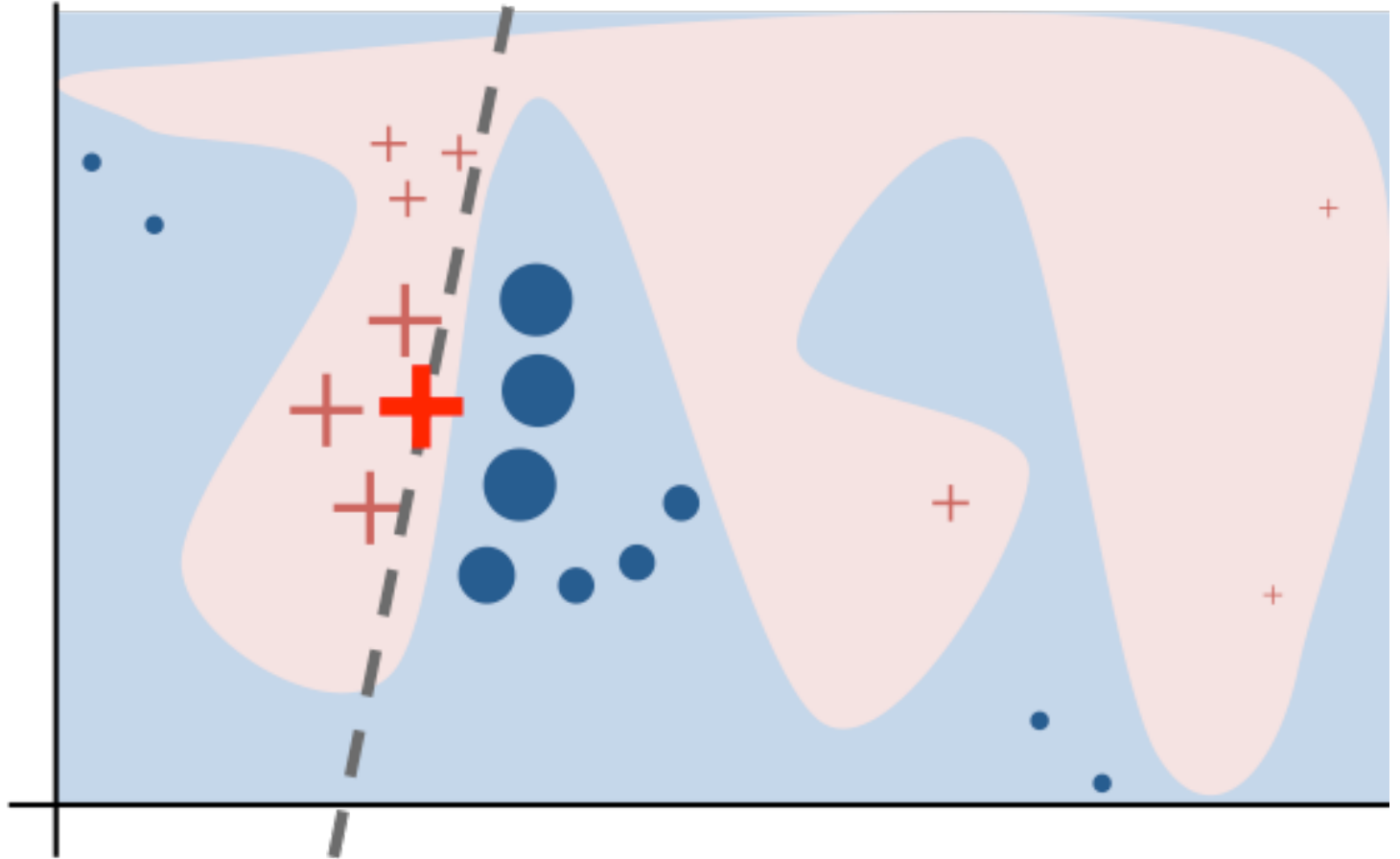
---

- Explinator: SM
  - Black Box: DNN, NN
  - Data Type: IMG
- Explinator: FI
  - Black Box: DNN, SVM
  - Data Type: ANY
- Explinator: DT
  - Black Box: ANY
  - Data Type: TAB

$R_1$ : IF(Outlook = Sunny) AND  
(Windy= False) THEN Play=Yes

# Local Explanation

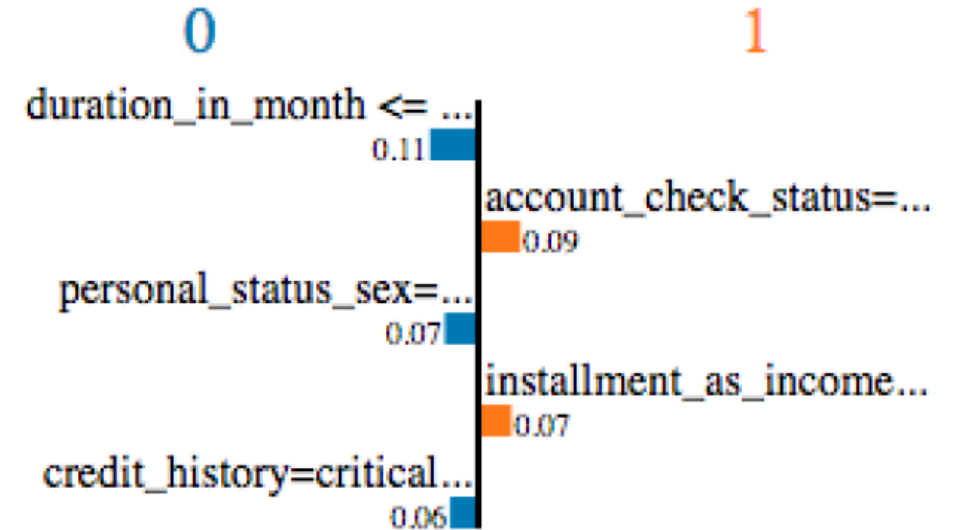
- The overall decision boundary is complex
- In the neighborhood of a single decision, the boundary is simple
- A single decision can be explained by auditing the black box around the given instance and learning a *local* decision.



# LIME – FI, AGN, ANY

```
01  Z = {}
02  x instance to explain
03  x' = real2interpretable(x)
04  for i in {1, 2, ..., N}
05      zi = sample_around(x')
06      z = interpretabel2real(z')
07      Z = Z ∪ {<zi, b(zi), d(x, z)>}
08  w = solve_Lasso(Z, k)
09  return w
```

black box  
auditing



- Marco Tulio Ribeiro, Sameer Singh, and Carlos Guestrin. 2016. Why should i trust you?: Explaining the predictions of any classifier. KDD.

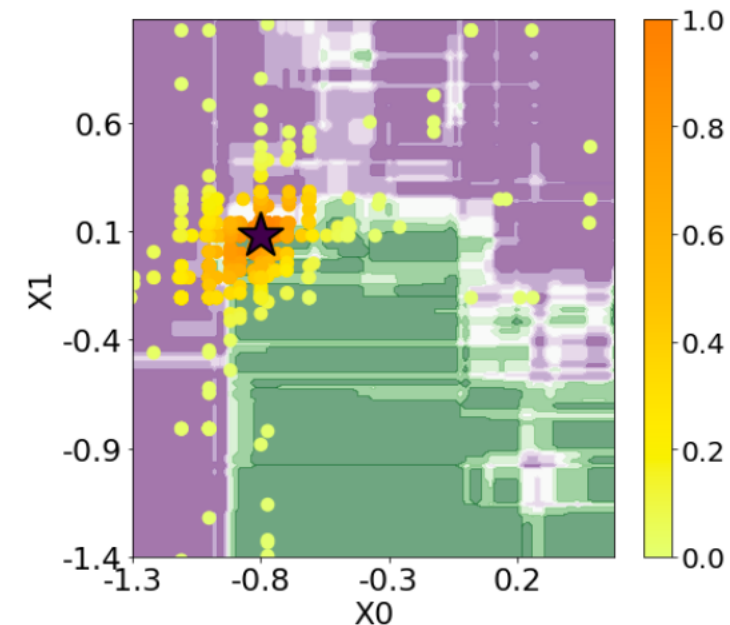
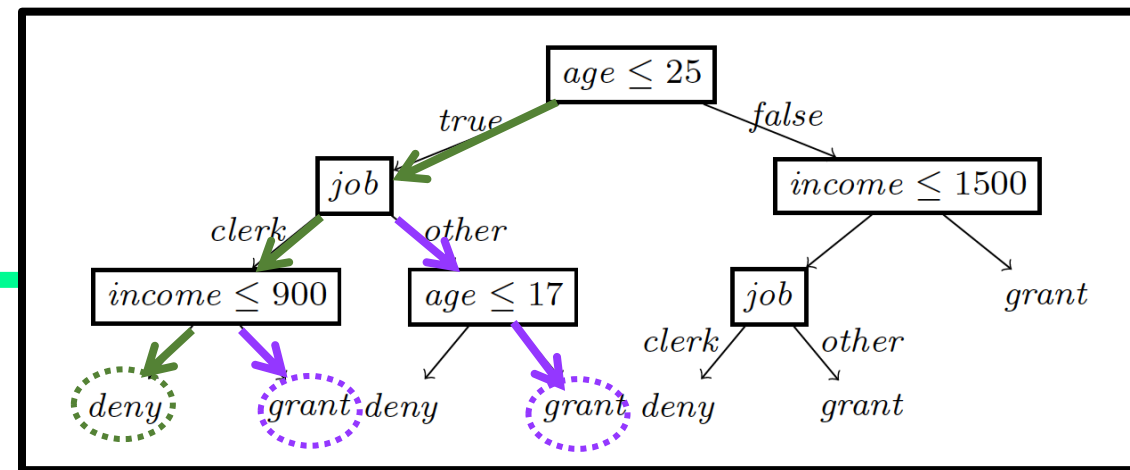
# LORE – DR, AGN, TAB

```
01  x instance to explain
02  Z= = geneticNeighborhood(x, fitness=, N/2)
03  Z≠ = geneticNeighborhood(x, fitness≠, N/2)
04  Z = Z= ∪ Z≠
05  c = buildTree(Z, b(Z)) black box auditing
06  r = (p -> y) = extractRule(c, x)
07  φ = extractCounterfactual(c, r, x)
08  return e = <r, φ>
```

$r = \{\text{age} \leq 25, \text{job} = \text{clerk}, \text{income} \leq 900\} \rightarrow \text{deny}$

$\Phi = \{(\{\text{income} > 900\} \rightarrow \text{grant}),$   
 $(\{17 \leq \text{age} < 25, \text{job} = \text{other}\} \rightarrow \text{grant})\}$

Pedreschi, Franco Turini,  
of *black box decision*





# Meaningful Perturbations – SM, DNN, IMG

- 01 `x` instance to explain
- 02 **varying** `x` into `x'` maximizing  $b(x) \sim b(x')$  ← *black box auditing*
- 03 the variation runs replacing a region `R` of `x` with:  
*constant value, noise, blurred image*
- 04 reformulation: find **smallest** `R` such that  $b(x_R) \ll b(x)$

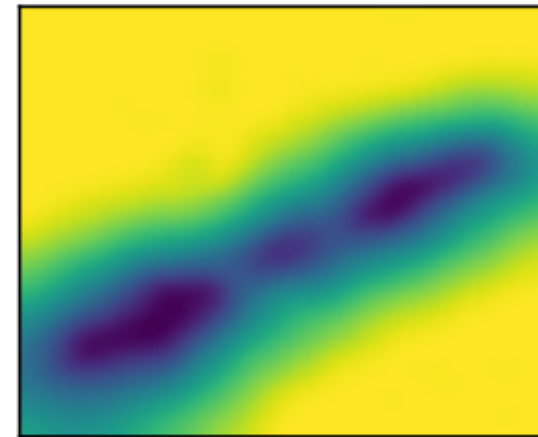
flute: 0.9973



flute: 0.0007

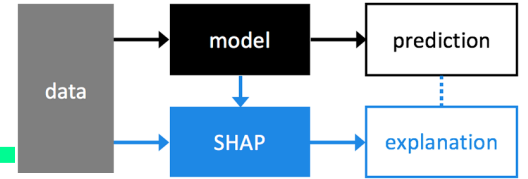


Learned Mask



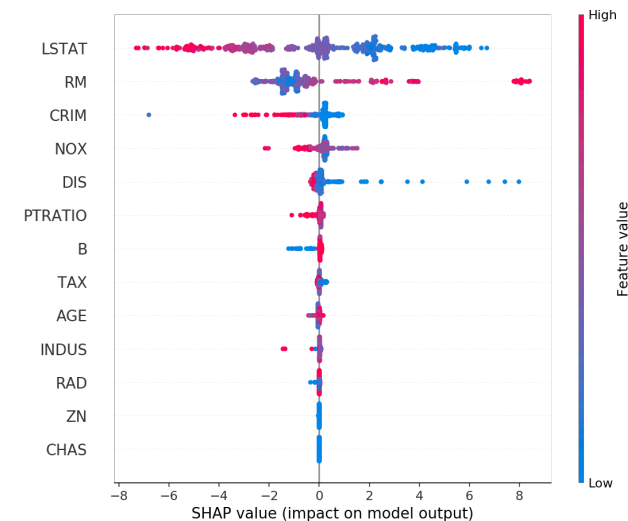
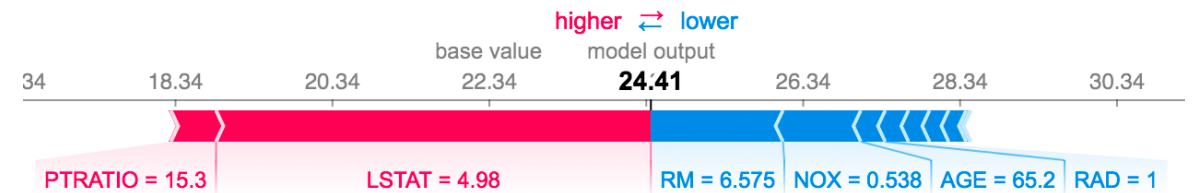


# SHAP (SHapley Additive exPlanations)



- SHAP assigns each feature an importance value for a particular prediction by means of an additive feature attribution method.
- It assigns an importance value to each feature that represents the effect on the model prediction of including that feature

$$g(z') = \phi_0 + \sum_{i=1}^M \phi_i z'_i,$$
$$\phi_i = \sum_{S \subseteq F \setminus \{i\}} \frac{|S|!(|F| - |S| - 1)!}{|F|!} [f_{S \cup \{i\}}(x_{S \cup \{i\}}) - f_S(x_S)]$$



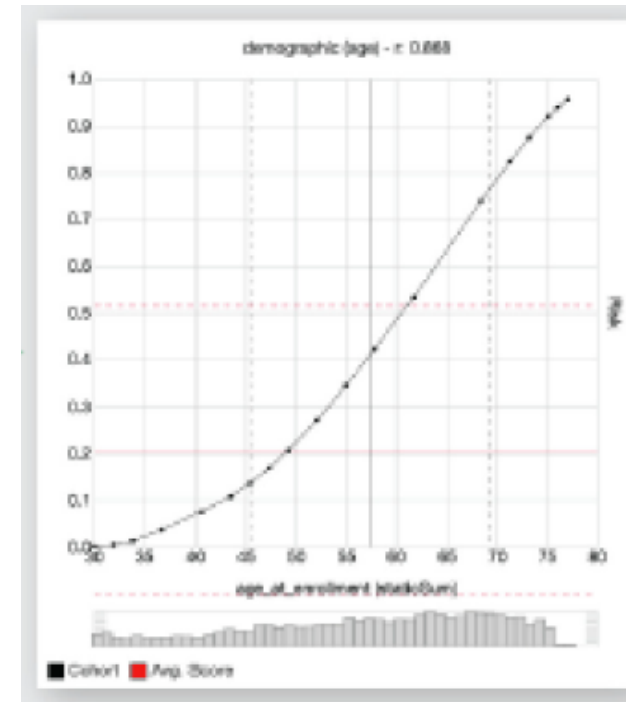
- Lundberg, Scott M., and Su-In Lee. "A unified approach to interpreting model predictions." *Advances in Neural Information Processing Systems*. 2017.

<i>Name</i>	<i>Ref.</i>	<i>Authors</i>	<i>Year</i>	<i>Explanator</i>	<i>Black Box</i>	<i>Data Type</i>	<i>General</i>	<i>Random</i>	<i>Examples</i>	<i>Code</i>	<i>Dataset</i>
NID	[83]	Olden et al.	2002	SA	NN	TAB			✓		
GDP	[8]	Baehrens	2010	SA	AGN	TAB	✓		✓		✓
QII	[24]	Datta et al.	2016	SA	AGN	TAB	✓		✓		✓
IG	[115]	Sundararajan	2017	SA	DNN	ANY			✓		✓
VEC	[18]	Cortez et al.	2011	SA	AGN	TAB	✓		✓		✓
VIN	[42]	Hooker	2004	PDP	AGN	TAB	✓		✓		✓
ICE	[35]	Goldstein et al.	2015	PDP	AGN	TAB	✓		✓	✓	✓
Prospector	[55]	Krause et al.	2016	PDP	AGN	TAB	✓		✓		✓
Auditing	[2]	Adler et al.	2016	PDP	AGN	TAB	✓		✓	✓	✓
OPIA	[1]	Adebayo et al.	2016	PDP	AGN	TAB	✓		✓		
—	[136]	Yosinski et al.	2015	AM	DNN	IMG			✓		✓
IP	[108]	Shwartz et al.	2017	AM	DNN	TAB			✓		
—	[137]	Zeiler et al.	2014	AM	DNN	IMG		✓		✓	
—	[112]	Springenberg et al.	2014	AM	DNN	IMG			✓		✓
DGN-AM	[80]	Nguyen et al.	2016	AM	DNN	IMG			✓	✓	✓

Solving The Model Inspection Problem

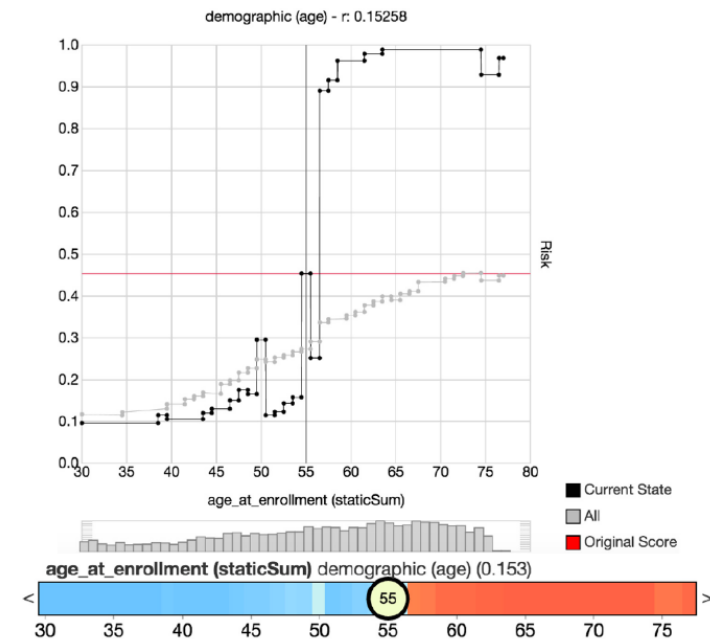
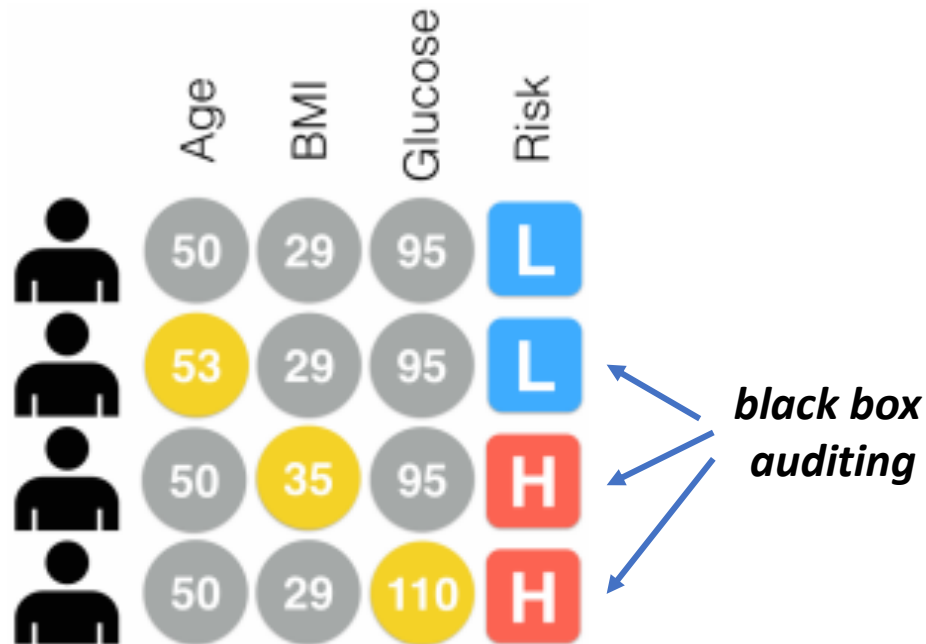
# Inspection Model Explainers

- Explinator: SA
  - Black Box: NN, DNN, AGN
  - Data Type: TAB
- Explinator: PDP
  - Black Box: AGN
  - Data Type: TAB
- Explinator: AM
  - Black Box: DNN
  - Data Type: IMG, TXT



# Prospector – PDP, AGN, TAB

- Introduce *random perturbations* on input values to understand to which extent every feature impact the prediction using PDPs.
- The input is changed *one variable at a time*.



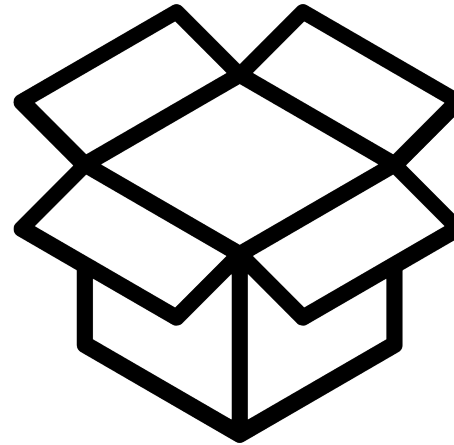
<i>Name</i>	<i>Ref.</i>	<i>Authors</i>	<i>Year</i>	<i>Explanator</i>	<i>Black Box</i>	<i>Data Type</i>	<i>General</i>	<i>Random</i>	<i>Examples</i>	<i>Code</i>	<i>Dataset</i>
CPAR	[135]	Yin et al.	2003	DR	—	TAB					✓
FRL	[127]	Wang et al.	2015	DR	—	TAB			✓	✓	✓
BRL	[66]	Letham et al.	2015	DR	—	TAB			✓		
TLBR	[114]	Su et al.	2015	DR	—	TAB			✓		✓
IDS	[61]	Lakkaraju et al.	2016	DR	—	TAB			✓		
Rule Set	[130]	Wang et al.	2016	DR	—	TAB			✓	✓	✓
1Rule	[75]	Malioutov et al.	2017	DR	—	TAB			✓		✓
PS	[9]	Bien et al.	2011	PS	—	ANY			✓		✓
BCM	[51]	Kim et al.	2014	PS	—	ANY			✓		✓
OT-SpAMs	[128]	Wang et al.	2015	DT	—	TAB			✓	✓	✓

Solving The Transparent Design Problem

# Transparent Model Explainers

---

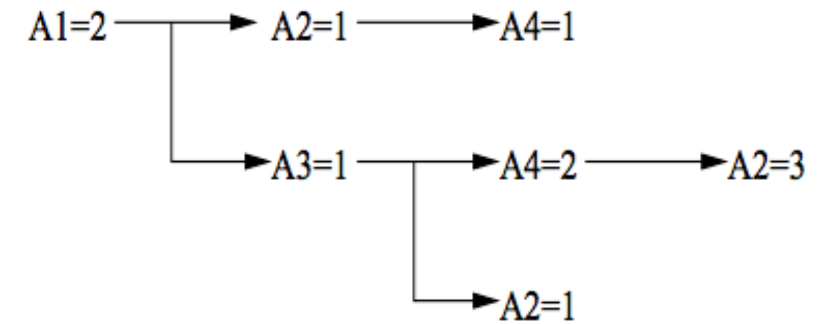
- Explanators:
  - DR
  - DT
  - PS
- Data Type:
  - TAB



# CPAR – DR, TAB

- Combines the advantages of associative classification and rule-based classification.
- It adopts a greedy algorithm to generate **rules directly from training data**.
- It generates more rules than traditional rule-based classifiers to **avoid missing important rules**.
- To **avoid overfitting** it uses expected accuracy to evaluate each rule and uses the best  $k$  rules in prediction.

$(A_1 = 2, A_2 = 1, A_4 = 1).$   
 $(A_1 = 2, A_3 = 1, A_4 = 2, A_2 = 3).$   
 $(A_1 = 2, A_3 = 1, A_2 = 1).$



# CORELS – DR, TAB

- It is a ***branch-and bound algorithm*** that provides the optimal solution according to the training objective with a certificate of optimality.
- It ***maintains a lower bound*** on the minimum value of error that each incomplete rule list can achieve. This allows to ***prune an incomplete rule list*** and every possible extension.
- It terminates with the optimal rule list and a certificate of optimality.

```
if (age = 18 – 20) and (sex = male) then predict yes
else if (age = 21 – 23) and (priors = 2 – 3) then predict yes
else if (priors > 3) then predict yes
else predict no
```



OPENING THE

Take Home Message

BLACK  
BOX

# Take-Home Messages

---

- Explainable AI is motivated by real-world application of AI
- Not a new problem – a reformulation of past research challenges in AI
- Multi-disciplinary: multiple AI fields, HCI, social sciences (multiple definitions)
- In Machine Learning:
  - Transparent design or post-hoc explanation?
  - Background knowledge matters!
  - We can scale-up symbolic reasoning by coupling it with representation learning on graphs.
- In AI (in general): many interesting / complementary approaches

# Open The Black Box!

---

- ***To empower*** individual against undesired effects of automated decision making
- ***To reveal*** and protect new vulnerabilities
- ***To implement*** the “right of explanation”
- ***To improve*** industrial standards for developing AI-powered products, increasing the trust of companies and consumers
- ***To help*** people make better decisions
- ***To align*** algorithms with human values
- ***To preserve*** (and expand) human autonomy



# Open Research Questions

---

- There is ***no agreement*** on ***what an explanation is***
- There is ***not a formalism*** for ***explanations***
- There is ***no work*** that seriously addresses the problem of ***quantifying*** the grade of ***comprehensibility*** of an explanation for humans
- Is it possible to join ***local*** explanations to build a ***globally*** interpretable model?
- What happens when black box make decision in presence of ***latent features***?
- What if there is a ***cost*** for querying a black box?



# References

---

- Guidotti, R., Monreale, A., Ruggieri, S., Turini, F., Giannotti, F., & Pedreschi, D. (2018). ***A survey of methods for explaining black box models***. *ACM Computing Surveys (CSUR)*, 51(5), 93
- Finale Doshi-Velez and Been Kim. 2017. ***Towards a rigorous science of interpretable machine learning***. arXiv:1702.08608v2
- Alex A. Freitas. 2014. ***Comprehensible classification models: A position paper***. ACM SIGKDD Explor. Newslett.
- Andrea Romei and Salvatore Ruggieri. 2014. ***A multidisciplinary survey on discrimination analysis***. Knowl. Eng.
- Yousra Abdul Alsahib S. Aldeen, Mazleena Salleh, and Mohammad Abdur Razzaque. 2015. ***A comprehensive review on privacy preserving data mining***. SpringerPlus
- Marco Tulio Ribeiro, Sameer Singh, and Carlos Guestrin. 2016. ***Why should i trust you?: Explaining the predictions of any classifier***. KDD.
- Houtao Deng. 2014. ***Interpreting tree ensembles with intrees***. arXiv preprint arXiv:1408.5456.
- Mark Craven and JudeW. Shavlik. 1996. ***Extracting tree-structured representations of trained networks***. NIPS.

# References

---

- M. Gethsiyal Augasta and T. Kathirvalavakumar. 2012. ***Reverse engineering the neural networks for rule extraction in classification problems***. NPL
- Riccardo Guidotti, Anna Monreale, Salvatore Ruggieri, Dino Pedreschi, Franco Turini, and Fosca Giannotti. 2018. ***Local rule-based explanations of black box decision systems***. arXiv preprint arXiv:1805.10820
- Ruth Fong and Andrea Vedaldi. 2017. ***Interpretable explanations of black boxes by meaningful perturbation***. arXiv:1704.03296 (2017).
- Paulo Cortez and Mark J. Embrechts. 2011. ***Opening black box data mining models using sensitivity analysis***. CIDM.
- Ruth Fong and Andrea Vedaldi. 2017. ***Interpretable explanations of black boxes by meaningful perturbation***. arXiv:1704.03296 (2017).
- Xiaoxin Yin and Jiawei Han. 2003. ***CPAR: Classification based on predictive association rules***. SIAM, 331–335
- Angelino, E., Larus-Stone, N., Alabi, D., Seltzer, M., & Rudin, C. 2017. ***Learning certifiably optimal rule lists***. KDD.